

This correspondence is being deposited with the United States Postal Service as Express Mail addressed to: Box Patent Application, Assistant Commissioner for Patents, Washington, DC 20231, on March 30, 2001, Express Mail Receipt No. EF 055069970 US.

1 MULTI-SERVICE SEGMENTATION AND REASSEMBLY DEVICE

2
3 Bidyut Parruck

4 Chulanur Ramakrishnan

5
6 CROSS-REFERENCE TO RELATED APPLICATION

7 This application claims the benefit under 35 U.S.C. §120 of U.S. Patent
8 Application Serial No. 09/779,381, filed February 7, 2001. The subject matter of
9 U.S. Patent Application Serial No. 09/779,381 is incorporated herein by reference.

10
11 TECHNICAL FIELD

12 This invention relates to multi-service network communications, such as, for
13 example, line card circuitry disposed within routers and switches.

14
15 BACKGROUND INFORMATION

16 Figure 1 (Prior Art) is a diagram of a part of the Internet 1. The Internet, loosely
17 defined, is a collection of networks that are interconnected by devices called
18 “routers”. In the illustration, the Internet 1 involves seven networks N1-N7 and five
19 routers R1-R5. A protocol called the Internet Protocol (IP) is used to communicate a
20 message from a source device (a node) on one network to a destination device (a
21 node) on another network. The message is broken up into pieces and each of these
22 pieces is packaged into what is called an “IP packet”. These packets may be of
23 varying lengths. The IP packets of the message are then sent from the source to the
24 destination from one network to the next via the routers. The various IP packets can
25 take different paths to get from the source to the destination. When all the IP
26 packets arrive at the destination, they are reassembled to recreate the original
27 message.

28 This high level IP message can be transported across an individual network using
29 any one of many lower level protocols. Some of the protocols are packet-based

1 protocols, whereas others of the protocols are cell-based protocols. One packet-
2 based protocol used to transport IP is called Multi-Protocol Label Switching (MPLS).
3 In MPLS, each packet is encapsulated with an MPLS label by the first MPLS device
4 it encounters as it enters an MPLS network. The MPLS device is called an MPLS
5 edge router. The MPLS edge router analyses the contents of the IP header and
6 selects an appropriate MPLS label with which to encapsulate the packet. MPLS
7 packets therefore have varying lengths in the same way that IP packets do. At all
8 the nodes within the network subsequent to the edge router, the MPLS label (and
9 not the IP header) is used to make the forwarding decisions for the packet. Paths
10 through particular nodes in the network are setup from edge to edge, the label
11 defining the particular path its packet will take. Finally, as an MPLS labeled packet
12 leaves the network via an edge router, the edge router removes the MPLS label.

13 One cell-based lower level protocol used to transport IP over a network is the
14 Asynchronous Transfer Mode (ATM) protocol. In ATM, all packets are of equal
15 length. They are therefore called "cells". A large IP packet is transported over an
16 ATM network by segmenting the large IP packet into a plurality of smaller pieces.
17 Each of the smaller pieces is packaged to become an ATM cell. The ATM cells are
18 then transported across the ATM network. When the ATM cells reach the edge of
19 the ATM network, their payloads are reassembled to reform the large IP packet. In
20 Figure 1, networks N1, N5 and N3 are cell-based ATM networks. Networks N2, N6,
21 N4 and N7 are packet-based MPLS networks.

22 In the example of Figure 1, networks N3 and N4 are OC-192 high-speed
23 networks adapted to carry traffic over long distances. Router R2 at one end of
24 network N3 may, for example, be located in San Francisco whereas router R4 at the
25 other end of network N3 may be located in New York. Such high-speed long
26 distance networks are often called the "backbone" of the Internet.

27 In the example of Figure 1, individual users U1-U10 are coupled to the Internet
28 via local area networks. Networks N1, N2 and N7 are local area networks. In one
29 example where the network is a corporate network serving an office building, the
30 users are corporate employees in a building. In an example where the network is a
31 network operated by an Internet Service Provider (ISP), the users are individual
32 customers that pay the ISP to gain access to the Internet.

1 Consider the situation where users on networks N1 and N2 issue IP messages
2 that are destined to go to destinations to the right side of the backbone such that the
3 messages should go through one of the two back bone networks N3 and N4. In
4 such a case, the IP traffic from networks N1 and N2 is aggregated and supplied to
5 the router access point on the appropriate backbone network. A portion of the
6 Internet called the "Metropolitan Area" performs this function. In the illustration, the
7 metro area includes a router R1 used for aggregating traffic from networks N1 and
8 N2, and for routing that information to the appropriate one of backbone networks N3
9 and N4.

10 Figure 2 (Prior Art) is a more detailed view of router R1. Router R1 includes line
11 cards 2-3 for interfacing to ATM networks, other line cards 4 and 5 for interfacing to
12 MPLS networks, and a switch fabric 11. ATM line card 3 is coupled to ATM network
13 N5 such that router R1 can communicate with backbone network N3 via network N5.
14 Similarly, MPLS line card 5 is coupled to MPLS network N6 such that router R1 can
15 communicate with backbone network N4 via network N6. ATM line card 2 is coupled
16 to ATM network N1 via OC-12 fiber optic link 6, SONET multiplexer 7, higher speed
17 OC-48 fiber optic link 8, and SONET multiplexer 9. MPLS line card 4 is coupled to
18 MPLS network N2 via OC-12 fiber optic link 10, SONET multiplexer 7, higher speed
19 OC-48 fiber optic link 8, and SONET multiplexer 9. SONET multiplexer 7 performs
20 time division multiplexing (TDM) to modulate both ATM traffic from network N1 as
21 well as packet MPLS traffic from network N2 onto the same wavelength channel
22 transmitted over the same fiber optic link 8. SONET multiplexer 9 performs the
23 inverse function of time demultiplexing the signal on fiber optic link 8 to extract and
24 separate the ATM traffic from the MPLS traffic.

25 Router R1, when it receives an IP message from one of networks N1 or N2,
26 determines whether to forward the message on the message's "next hop" to router
27 R2 or R3. In this way IP network information from the users is aggregated in the
28 metro area and is directed to the correct backbone network for communication
29 across long distances to the destination.

30 A problem may exist if one of the local area networks coupled to router R1 is
31 disconnected or if the type of traffic on that network is changed from MPLS packet
32 traffic to ATM cell traffic or visa versa. Consider the situation in which ATM network

1 N1 ceases business operations. In that case, the operator of router R1 will likely
2 want to disconnect network N1 from its SONET multiplexer 7 and to couple in the
3 network of another paying customer. For example, the operator may want to
4 disconnect ATM network N1 and to connect in its place MPLS network N7. If this is
5 done, however, MPLS traffic would be received on ATM line card 2. ATM line card 2
6 is not suitable for coupling to an MPLS network. Consequently, ATM line card 2
7 may have to be disconnected and a suitable MPLS line card substituted in its place.
8 With the expansion of the Internet and with advances in IP switching technology, it
9 appears that the proportion of packet networks to ATM networks is increasing.
10 Accordingly, as more and more of the networks coupled to a router such as router
11 R1 migrate from one type of traffic to the other, more and more of the line cards of
12 the router will have to be replaced. This is undesirable. A solution is desired whereby a
13 smooth and easy migration from one type of traffic to the next is possible without the
14 removal of line cards or the physical manipulation of the router.

15 Figure 3 is a diagram of one possible approach to the problem involving a line
16 card 12 that handles both ATM and packet traffic. Line card 12 is coupled to a
17 switch fabric of a router by interface 13. Cell and packet traffic received from fiber
18 optic cable 14 and transmitted on fiber optic cable 15 are time division
19 multiplexed/demultiplexed by TDM device 16. Cell traffic is handled by integrated
20 circuit device 17. Packet traffic is handled by integrated circuit device 18. As the
21 relative amounts of cell traffic to packet traffic change, the same line card can be
22 used.

23 24 SUMMARY

25 Although the line card set forth in Figure 3 is satisfactory for some applications,
26 the general approach of Figure 3 involves substantial cost. The non-recurring
27 engineering costs associated with developing an integrated circuit can be quite high.
28 Even though the development of a particular integrated circuit may make technical
29 sense for a given application, it may be economically unreasonable to do so where
30 production volumes of the integrated circuit would be low. For each line card of
31 Figure 3, there is one integrated circuit for handling ATM traffic and one integrated
32 circuit for handling packet traffic. Developing a single integrated circuit having the

1 functionality of both the ATM device and the packet device might involve less non-
2 recurring engineering costs than developing two different integrated circuits, but the
3 integration would likely result in an undesirably large integrated circuit. Parts of such
4 a single integrated circuit may see little use in certain circumstances. Consider the
5 situation in which the mix of ATM traffic to packet traffic shifts to where there is little
6 or no ATM traffic. The single integrated circuit would involve a data path and
7 associated circuitry for handling ATM traffic that is underutilized or is not used at all.
8 Providing this extra unnecessary circuitry on the single integrated circuit would
9 constitute an unnecessary expense. It would therefore make the line card solution
10 afforded by the single integrated circuit unnecessarily expensive. It would be
11 preferable to get more use out of the circuitry provided on the integrated circuit in
12 order to reduce costs.

13 Not only might some circuitry be underutilized, but also so might other circuitry
14 become overburdened. In the above situation, for example, more and more
15 processing responsibilities would be required of the packet handling circuitry. If the
16 processing capability of the packet handling circuitry is sized to accommodate a
17 particular load, then increasing the load beyond that load by more and more of the
18 traffic shifting to packet traffic may result in the line card being inadequate to handle
19 the required traffic.

20 In one novel aspect, the same circuitry on a single Multi-Service Segmentation
21 And Reassembly (MS-SAR) integrated circuit handles both cell traffic and packet
22 traffic. Rather than there being a first data path through the integrated circuit for cell
23 processing, and another data path through the integrated circuit for packet
24 processing, functional blocks along a single data path process cell and/or packet
25 information that passes through the integrated circuit through the single data path.
26 Individual flows passing through the single data path are processed in accordance
27 with the traffic type of the individual flow. Any mix of cell to packet traffic can be
28 accommodated, thereby enabling a smooth migration from one traffic type to
29 another. The MS-SAR can handle dynamic changes in the relative amounts of cell
30 and packet traffic. Production costs associated with the integrated circuit are
31 reduced because the same functional blocks are used, albeit in different ways, to
32 process different types of flows.

1 In another novel aspect, production volumes may be increased in order to realize
2 economies of scale and to reduce per part cost. To allow production volumes to be
3 increased, the very same MS-SAR integrated circuit is usable not only both as an
4 ingress device and an egress device, but also with either a packet-based switch
5 fabric or a cell-based switch fabric. By providing a single MS-SAR integrated circuit
6 that functions as either an ingress device or an egress device as compared to a
7 device that can function as just one or the other, the number of applications for the
8 MS-SAR integrated circuit is increased. By providing a single MS-SAR integrated
9 circuit that can work with either a packet-based switch fabric or a cell-based switch
10 fabric, the number of applications for the single MS-SAR integrated circuit is
11 increased as compared to the same device that could only function with one type of
12 switch fabric.

13 In another novel aspect, a single MS-SAR integrated circuit involves a lookup
14 block, a segmentation block, and a reassembly block, wherein traffic flows through a
15 single data path through the lookup block, and then through the segmentation block,
16 and then through the reassembly block. By using two identical such MS-SAR
17 devices, one in an ingress mode and another in an egress mode, information
18 received onto the ingress MS-SAR in either ATM or packet format can be
19 communicated through either a packet-based or a cell-based switch fabric, and can
20 be output from the egress MS-SAR in either ATM or packet format. Information
21 communicated in AAL5 adaptation cells received onto the ingress MS-SAR can be
22 reassembled and output in packet format on the egress MS-SAR. Packet format
23 information received onto the ingress MS-SAR can be segmented and output in the
24 form of AAL5 adaptation cells from the egress MS-SAR. The versatility of the single
25 MS-SAR to handle many different traffic types using the same circuitry further
26 increases the number of applications for which the integrated circuit can be used.

27 In another novel aspect, individual flows are processed in different ways by an
28 egress MS-SAR. An indication of the type of egress processing to be done by an
29 MS-SAR on a flow is embedded into the flow that is received from a switch fabric
30 onto the egress MS-SAR. The egress MS-SAR reads the embedded indication and
31 performs the type of egress processing indicated. In one embodiment, the indication
32 of the type of egress processing is contained in a switch header (the switch header

1 can be either a switch header for a cell-based switch fabric or for a packet-based
2 switch fabric), the switch header being added by a first MS-SAR functioning in an
3 ingress mode, the first MS-SAR and the second MS-SAR being substantially
4 identical integrated circuits. In one embodiment, information on how to locate the
5 indication of the type in the information coming into the egress MS-SAR is provided
6 to the egress MS-SAR for each logical input port of the egress MS-SAR. The egress
7 MS-SAR uses: 1) the input port number of the flow, and 2) the information on how to
8 locate the indication for a given input port, to locate the indication in the information
9 coming in from the switch fabric.

10 In accordance with known ATM adaptation layer protocols, packet data can be
11 transmitted over an ATM network by segmenting the packet into smaller pieces and
12 then transmitting each smaller piece over the ATM network in the form of an ATM
13 cell. After transmission across the ATM network, the data payloads of the individual
14 ATM cells are recovered and are reassembled into the original packet. This
15 segmentation and reassembly process has been traditionally performed on a line
16 card by reassembling packets as the individual ATM cells are received onto the line
17 card. A reassembly context is maintained for each packet being reassembled. This
18 reassembly context may, for example, include a composite cyclic redundancy check
19 (CRC) value that is modified as each ATM cell carrying a portion of the packet is
20 received and processed. The CRC value calculated on the data portions received is
21 checked against a CRC transmitted in a trailer of the last ATM cell to verify that the
22 data carried by the numerous ATM cells has not been corrupted in the segmentation
23 and reassembly process. Once a packet has been checked and reassembled, is it
24 buffered into a payload memory on the line card. If, for example, the line card were
25 to support the simultaneous reassembly of one million packets, then the line card
26 would have to be able to store one million reassembly contexts. This would involve
27 a large amount of memory.

28 In another novel aspect, packets to be reassembled on the line card in such an
29 adaptation layer process are not reassembled before being buffered on the line card.
30 Rather, the payloads of the individual cells are buffered in payload memory as cells.
31 The line card does not maintain a reassembly context for each such packet being
32 buffered. When the buffered packet information is to be output from the line card,

1 the cell payloads corresponding to the packet are read out of payload memory and
2 the cell payloads are reassembled to form the packet. In this way, a large number of
3 packets being received onto the line card do not need to be simultaneously
4 reassembled, but rather the number of packets being simultaneously reassembled
5 can be set to have a smaller maximum. In one embodiment, one million flows to be
6 reassembled can be received at one time onto a line card, but only one packet per
7 active line card output port is reassembled at a time. By reducing the maximum
8 number of packets being simultaneously reassembled, the maximum number of
9 reassembly contexts to be stored on the line card is reduced. Reducing the number
10 of reassembly contexts to be stored reduces the amount of memory necessary and
11 thereby reduces line card costs.

12 Not only is a reassembly context involved the reassembly process, but a
13 segmentation context is also traditionally involved in the segmentation process.
14 Traditionally, packets to be segmented on a line card in accordance with an
15 adaptation layer process are received and stored into payload memory as packets.
16 When a packet is to be output, it is retrieved from payload memory and is
17 segmented into segments such that each of the segments forms the data payload of
18 an ATM cell. To be able to check the integrity of the data when the segments are
19 reassembled, a CRC is calculated on the packet at the time of segmentation and is
20 transmitted in a trailer that is included in the last ATM cell. For each such
21 segmentation process that is going on at the same time, a segmentation context
22 including a partial CRC value is maintained. If a large number such as a million
23 simultaneous output flows is to be supported by the line card, then the large number
24 of segmentation contexts must be stored on the line card. This involves a lot of
25 memory and consequently increases line card cost.

26 In another novel aspect, packets to be segmented in accordance with an
27 adaptation layer protocol are segmented on a per input port basis as they are
28 received onto the line card and prior to being buffered on the line card. The packets
29 are not buffered on the line card, but rather segments are buffered. Because only
30 one segmentation is performed at a time for a given line card input port, the
31 maximum number of simultaneous segmentations is limited to the number of input
32 ports. By limiting the maximum number of simultaneous segmentations to be

1 performed, the memory required to store the associated segmentation contexts is
2 reduced. In one embodiment, one million simultaneous flows can be processed, but
3 there are only sixty-four input ports. The amount of memory required for storing
4 segmentation contexts is therefore significantly reduced.

5 In another novel aspect, an MS-SAR involves a data path such that data received
6 onto the MS-SAR passes through the data path and to a memory manager that
7 stores the data into payload memory. The data is then read out of the payload
8 memory and passes through the remainder of the data path to be output from the
9 MS-SAR. How and when the data is read out of payload memory is controlled by
10 control circuitry. The control circuitry controls the memory manager so that the
11 memory manager retrieves data from payload memory so that it will be output from
12 the MS-SAR in a manner controlled by the control circuitry. In one novel aspect, the
13 MS-SAR is partitioned into two integrated circuits such that the data path circuitry is
14 disposed on one integrated circuit and such that the control circuitry is disposed on
15 another integrated circuit. Partitioning the MS-SAR in this way facilitates future
16 increasing of data throughput rates without redesigning the MS-SAR. To increase
17 data throughput, for example from OC-192 rates to OC-768 rates, multiple data path
18 integrated circuits are disposed in parallel, each of the data path integrated circuits
19 being controlled by the same control integrated circuit. The control integrated circuit
20 has multiple control interfaces, one such interface for coupling to and controlling
21 each of the data path integrated circuits.

22 In another novel aspect, a router involves a first line card and a second line card.
23 Each of the first and second line cards involves an MS-SAR operating in the ingress
24 mode and an MS-SAR operating in the egress mode. The egress MS-SAR on the
25 second line card can become endangered of being overloaded if, for example, the
26 ingress MS-SAR on the first line card continues to send network information for a
27 flow to the egress MS-SAR on the second line card, but the egress MS-SAR on the
28 second line card is prevented from outputting that information, for example due to
29 congestion at the framer. Consequently, more and more of the network information
30 for the flow will end up having to be buffered by the egress MS-SAR of the second
31 line card. In one novel aspect, the ingress and egress MS-SAR devices of the first
32 line card are linked by a serial bus on the first line card, and the ingress and egress

1 MS-SAR devices of the second line card are linked by a serial bus on the second
2 line card. If the egress MS-SAR of the second line card is in danger of becoming
3 overloaded, then the egress MS-SAR of the second line card sends an indication of
4 this situation to the ingress MS-SAR of the second line card via the serial bus on the
5 second line card. The ingress MS-SAR of the second line card relays that
6 information to the first line card by outputting a special status switch cell. The
7 special status switch cell is transported across the switch fabric to the egress MS-
8 SAR of the first line card. The egress MS-SAR of the first line card detects the
9 special status switch cell, and relays the indication of the situation to the ingress MS-
10 SAR of the first line card via the serial bus on the first line card. In response to
11 receiving this indication from the serial bus, the ingress MS-SAR on the first line card
12 slows or stops outputting the information that is overburdening the egress MS-SAR
13 on the second line card.

14 In another novel aspect, an integrated circuit includes a reassembly circuit,
15 wherein in an ingress mode the reassembly circuit uses a flow ID to lookup a switch
16 header in an memory, and wherein in an egress mode the reassembly circuit uses a
17 flow ID to lookup a network protocol header (for example, ATM header or MPLS
18 header) in the memory.

19 These are but some of the novel aspects. Other novel structures and methods
20 are disclosed in the detailed description below. This summary does not purport to
21 define the invention. The invention is defined by the claims.

22 BRIEF DESCRIPTION OF THE DRAWINGS

23 Figure 1 (Prior Art) illustrates a part of the Internet where different traffic types are
24 aggregated.

25 Figure 2 (Prior Art) is a diagram of a router in the part of the Internet illustrated in
26 Figure 1.

27 Figure 3 is a diagram illustrative of an approach to solving a problem associated
28 with the network structure of Figure 1. Although the line card as illustrated may and
29 is intended to reflect a product designed by or being designed by Northern Telecom
30 Ltd., adequate details about this line card are not available to the inventors or the
31

1 assignee to state here in this patent document that what is shown in Figure 3 is prior
2 art.

3 Figure 4 is a diagram of a switching device (in this case a router) in accordance
4 with an embodiment of the present invention.

5 Figure 5 is a diagram of a line card in the router of Figure 4.

6 Figure 6 is a diagram that sets forth various application types that the router of
7 Figure 4 (involving a pair of MS-SAR devices, one operating in an ingress mode,
8 and the other operating in an egress mode) can carry out.

9 Figure 7 is a diagram illustrating ingress and egress application types involving a
10 cell-based switch fabric. These application types can be carried out by the router of
11 Figure 4.

12 Figure 8 is a diagram illustrating ingress and egress application types involving a
13 packet-based switch fabric. These application types can be carried out by the router
14 of Figure 4.

15 Figure 9 is a diagram of an example in accordance with an embodiment of the
16 present invention wherein a first flow is processed in accordance with ingress
17 application type 3 and egress application type 11, and wherein a second flow is
18 processed in accordance with ingress application type 0 and egress application type
19 8.

20 Figure 10 is a diagram illustrating an MS-SAR and associated memories in
21 accordance with an embodiment of the present invention.

22 Figure 11 is a simplified representation of an MPLS packet of flow #1 in the
23 example of Figure 9.

24 Figure 12 is a simplified diagram of chunks of flow #1 that are output by the
25 incoming SPI-4 interface block in the example of Figure 9.

26 Figure 13 is a diagram that illustrates information flow into an ingress MS-SAR in
27 the example of Figure 9.

28 Figure 14 is a representation of a port table in the lookup engine.

29 Figure 15A-15C are diagrams of the building of a per flow queue for flow #1 in the
30 example of Figure 9.

31 Figure 16 is a diagram of a dequeue memory location (stores the header pointer)
32 for one FID.

Figures 17 and 18 are diagrams of the first and second enqueue memory locations (store the tail pointer) for one FID.

Figure 19 is a diagram of a memory location (stores a pointer to a buffer in the body of a queue) for an intermediate BID in a linked list.

Figure 20 is a diagram of an ATM cell of flow #2 in the example of Figure 9.

Figure 21 is a diagram of a 56-byte chunk of flow #2 as output from the incoming SPI-4 interface block of the ingress MS-SAR in the example of Figure 9.

Figure 22 is a diagram of a 64-byte chunk of flow #2 as output from the segmentation block of the ingress MS-SAR in the example of Figure 9.

Figure 23 is a diagram of a per flow queue for flow #2 of the example of Figure 9.

Figure 24 is a diagram of the port calendar in the reassembly block in the MS-SAR.

Figures 25 and 26 are diagrams of the port empty and port full registers in the reassembly block in the MS-SAR.

Figure 27 is a diagram that illustrates information flow from payload memory out of the ingress MS-SAR in the example of Figure 9.

Figure 28 is a diagram that illustrates the three switch cells of flow #1 in the example of Figure 9.

Figure 29 is a diagram of the switch cell for flow #2.

Figure 30 is a diagram that illustrates the general flow of information into egress MS-SAR 200 in the example of Figure 9.

Figure 31 is a diagram of the first switch cell for flow #1 in the example of Figure 9.

Figure 32 is a diagram that illustrates the general flow of information out of egress MS-SAR 200 in the example of Figure 9.

Figure 33 is a diagram that illustrates the format of one FID entry in the header table of the MS-SAR.

Figure 34 illustrates three 64-byte chunks of flow #1 as the chunks pass from reassembly block 205 to outgoing SPI-4 interface block 206 in the example of Figure 9.

Figure 35 illustrates the MPLS packet of flow #1 as output from framer 142 in the example of Figure 9.

1 Figure 36 illustrates the ATM cell of flow #2 as output from reassembly block 205
2 in the example of Figure 9.

3 Figure 37 illustrates the ATM cell of flow #2 as output from outgoing SPI-4
4 interface block 206 of egress MS-SAR 200 in the example of Figure 9.

5 Figure 38 illustrates an example of application types 5, 6, 14 and 13.

6 Figure 39 illustrates the processing of flow #1 in accordance with ingress
7 application type 5 in the example of Figure 38.

8 Figure 40 illustrates the processing of flow #2 in accordance with ingress
9 application type 6 in the example of Figure 38.

10 Figure 41 illustrates the processing of flow #1 in accordance with egress
11 application type 14 in the example of Figure 38.

12 Figure 42 illustrates the processing of flow #2 in accordance with egress
13 application type 13 in the example of Figure 38.

14 Figure 43 illustrates an example wherein a flow is processed on an ingress MS-
15 SAR in accordance with ingress application type 1 and is processed on an egress
16 MS-SAR in accordance with egress application type 9.

17 Figure 44 illustrates an example wherein a flow is processed on an ingress MS-
18 SAR in accordance with ingress application type 2 and is processed on an egress
19 MS-SAR in accordance with egress application type 10.

20 Figure 45 illustrates an example wherein a flow is processed on an ingress MS-
21 SAR in accordance with ingress application type 4 and is processed on an egress
22 MS-SAR in accordance with egress application type 14.

23 Figure 46 illustrates an example wherein a flow is processed on an ingress MS-
24 SAR in accordance with ingress application type 6 and is processed on an egress
25 MS-SAR in accordance with egress application type 12.

26 Figure 47 illustrates an example wherein a flow is processed on an ingress MS-
27 SAR in accordance with ingress application type 6 and is processed on an egress
28 MS-SAR in accordance with egress application type 14.

29 Figure 48 is a diagram of an embodiment wherein MS-SAR functionality is
30 partitioned into two integrated circuit chips (a data path integrated circuit and a
31 control integrated circuit) such that multiple data path integrated circuits chips can be
32 used with one control integrated circuit chip to increase data path throughput.

1 Figure 49 is a diagram of a packet as it is output from the distribution integrated
2 circuit of Figure 48.

3 Figures 50 and 51 are diagrams that illustrate the building of a packet queue in
4 connection with the operation of the embodiment of Figure 48.

5 Figure 52 is a diagram that illustrates a technique for accessing certain
6 information stored in external memory in a reduced amount of time in connection
7 with the embodiment of Figure 48.

8 Figure 53 is a diagram that illustrates a serial bus that couples an egress MS-
9 SAR of a line card to an ingress MS-SAR of the same line card. The egress MS-
10 SAR can use the serial bus to backpressure the sending ingress MS-SAR.

11 12 DETAILED DESCRIPTION

13 Figure 4 is a simplified diagram of a router 100 in accordance with an
14 embodiment of the present invention. Router 100 includes a plurality of line cards
15 101-104, a switch fabric 105 and a central processing unit (CPU) 106. The line
16 cards 101-104 are coupled to switch fabric 105 by parallel buses 107-114. In the
17 present example, each of parallel buses 107-114 is a 16-bit SPI-4, Phase II, LVDS
18 parallel bus operating at 400 MHz at a double data rate (DDR). CPU 106 is coupled
19 to line cards 101-104 by another parallel bus 131. In the present example, parallel
20 bus 130 is a 32-bit PCI bus. In this example, each of the line cards can receive
21 network communications in multiple formats. For example, line card 101 is coupled
22 to a fiber optic cable 115 such that line card 101 can receive from cable 115 network
23 communications at OC-192 rates in packets, ATM cells, and/or AAL5 cells. AAL5
24 cells are considered a type of ATM cell.

25 Line card 101 is also coupled to a fiber optic cable 116 such that line card 101
26 can output onto cable 116 network communications at OC-192 rates in packets,
27 ATM cells, and/or AAL5 cells. The fiber optic cables 117 and 118 across which line
28 card 103 communicates are also labeled in the diagram. All the line cards 101-104
29 in this example have substantially identical circuitry.

30 Figure 5 is a more detailed diagram of representative line card 101. Line card
31 101 includes OC-192 optical transceiver modules 119 and 120, two serial-to-parallel
32 devices (SERDES) 121 and 122, a framer integrated circuit 123, a IP classification

1 engine 124, two multi-service segmentation and reassembly devices (MS-SAR
 2 devices) 125 and 126, static random access memories (SRAMs) 127 and 128, and
 3 dynamic random access memories (DRAMs) 129 and 130. MS-SAR devices 125
 4 and 126 are identical integrated circuit devices, one of which (MS-SAR 125) is
 5 configured to be in an “ingress mode”, the other of which (MS-SAR 126) is
 6 configured to be in an “egress mode”. Each MS-SAR device includes a mode
 7 register that is written to by CPU 106 via bus 131. When router 100 is configured,
 8 CPU 106 writes to the mode register in each of the MS-SAR devices on each of the
 9 line cards so as to configure the MS-SAR devices of the line cards appropriately.

10 11 ROUTER AND LINE CARD:

12 Fiber optic cable 115 of Figure 4 can carry information modulated onto one or
 13 more of many different wavelengths (sometimes called “colors”). Each wavelength
 14 can be thought of as constituting a different communication channel for the flow of
 15 information. Accordingly, optics module 119 converts optical signals modulated onto
 16 one of these wavelengths into analog electrical signals. Optics module 119 outputs
 17 the analog electrical signals in serial fashion onto a high-speed analog bus 132.
 18 Serdes 121 receives this serial information and outputs it in parallel form to framer
 19 123 via high-speed parallel bus 133. Framer 123 receives the information from bus
 20 133, frames it, and outputs it to classification engine 124 via another SPI-4 bus 134.
 21 Classification engine 124 performs IP classification and outputs the information to
 22 the ingress MS-SAR 125 via another SPI-4 bus 135. The ingress MS-SAR 125
 23 processes the network information in various novel ways (explained below), and
 24 outputs the network information to switch fabric 105 (see Fig. 4) via SPI-4 bus 107.
 25 All the SPI-4 buses of Figures 4 and 5 are separate SPI-4, phase II, 400 MHz DDR
 26 buses having sixteen bit wide data buses.

27 Switch fabric 105, once it receives the network information, supplies that
 28 information to one of the line cards of router 100. Each of the line cards is identified
 29 by a “virtual output port” number. Router 100 can include up to 256 line cards.
 30 Accordingly, the virtual output port number has a range of from 0 to 255.

31 In the example of Figure 4, switch fabric 105 may supply network information
 32 received on fiber optic cable 115 to any one of fiber optic cables 116, 118, 136 or

1 137. It is one of the primary functions of router 100 to determine, based on the
2 certain information such as the intended destination of the network information, how
3 to route the information. In the case where the network information is in the IP
4 packet format, the router makes its decision on where to route the network
5 information based on an intended IP destination address present in the IP header of
6 each packet. In the case where the network information is in the ATM cell format,
7 the router makes its decision on where to route the network information based on a
8 virtual path identifier and virtual channel identifier (VPI/VCI) information in the ATM
9 header of each ATM cell.

10 If, for example, the network information received from fiber optic cable 115 is to
11 be output from fiber optic cable 118, then switch fabric 105 would direct that network
12 information to the "virtual output port" of line card 103. If, on the other hand, the
13 network information received from fiber optic cable 110 is to be output from fiber
14 optic cable 137, then switch fabric 105 would direct that network information to the
15 "virtual output port" of line card 104. To facilitate the rapid forwarding of such
16 network information through the switch fabric 105, network information passed to the
17 switch fabric 105 for routing is provided with a "switch header". The "switch header"
18 may be in a format specific to the manufacturer of the switch fabric of the router.
19 The switch header identifies the "virtual output port" to which the associated network
20 information should be routed. Switch fabric 105 uses the virtual output port number
21 in the switch header to route the network information to the correct line card.

22 It is, however, generally the router 100 that determines to which of the multiple
23 line cards particular network information will be routed. Accordingly, the router's
24 CPU 106 provisions lookup information in (or accessible to) the ingress MS-SAR
25 125 so that the MS-SAR 125 will append an appropriate switch header onto the
26 network information before the network information is sent to the switch fabric 105
27 for routing. The ingress MS-SAR 125 uses header information in the network
28 information to lookup the particular switch header specified by CPU 106 for the
29 particular network information. Once the MS-SAR 125 has found the switch header
30 that it is to append to the network information, MS-SAR 125 appends this switch
31 header to the network information and sends the network information with the switch
32 header on to the switch fabric 105 via SPI-4 bus 107.

1 Switch fabric 105 receives the network information and forwards it to the line card
2 identified by the particular "virtual output port" in the switch header. Consider the
3 example in which the network information is forwarded to the virtual output port of
4 line card 103 so that the network information from fiber optic cable 115 will be output
5 onto fiber optic cable 117. Because the structures of the various line cards are
6 identical in this example, the flow of network information through line card 103 is
7 explained as passing through the egress MS-SAR 126 of Figure 5. The network
8 information and switch header is received onto the egress MS-SAR 126 of the line
9 card that is identified by the virtual output port number in the switch header. The
10 egress MS-SAR 126 receives the network information, removes the switch header,
11 performs other novel processing (explained below) on the network information, and
12 outputs the network information to framer 123 via SPI-4 bus 138. Framer 123
13 outputs the network information to serdes 122 via high-speed parallel bus 139.
14 Serdes 122 converts the network information into serial analog form and outputs it to
15 output optics module 120 via high-speed analog bus 140. Output optics module 120
16 converts the information into optical signals modulated onto one wavelength
17 channel. This optical information is transmitted through fiber optic cable 116 (which
18 corresponds to cable 118).

19 Although each line card in the example of Figure 4 outputs to only one fiber optic
20 cable (i.e., has only one physical output port), this need not be the case. If, for
21 example, a high speed OC-192 optical transceiver and cable is coupled to a line
22 card, then the line card may only have one physical output port. If, on the other
23 hand, two lower speed optical transceivers (for example, OC-48 transceivers) and
24 cables are coupled to the line card, then the line card may have two physical output
25 ports (one being the first OC-48 transceiver and cable, the other being for the
26 second OC-48 transceiver and cable). The egress MS-SAR 126 outputs onto SPI-4
27 bus 138 network information to be transmitted along with an output port number from
28 0 to 63. (This "output port number" is not to be confused with the "virtual output port
29 number" used to identify line cards.) Groups of these 64 output port numbers are
30 mapped by framer 123 to the physical output ports provided on the line card. For
31 example, output ports 0-31 on SPI-4 bus 138 may be mapped to a first physical
32 output port (a first OC-48 output transceiver), whereas output ports 32-63 on SPI-4

bus 138 may be mapped to a second physical output port (a second OC-48 output transceiver). In accordance with the SPI-4 bus protocol, there are 64 logical output ports that egress MS-SAR 126 can specify on SPI-4 bus 138.

In a similar way, a line card may receive network information from more than one optical input transceiver module and more than one fiber optic cable. Information coming into a first optical input transceiver module (i.e., a first physical port) would be supplied onto SPI-4 bus 134 by framer 123 with a logical input port number that maps to the first optical input transceiver module onto which the information was received. Similarly, information coming into a second optical input transceiver module (i.e., a second physical port) would be supplied onto SPI-4 bus 134 by framer 123 with a logical input port number that maps to the second optical input transceiver module onto which the information was received. In accordance with the SPI-4 bus protocol, there are 64 logical input ports that framer 123 can specify on SPI-4 bus 134. The logical input port information on SPI-4 bus 134 flows through classification engine 124 to appear on SPI-4 bus 135.

APPLICATION TYPES:

A flow is a sequence of one or more packets or cells transmitted between a selected source and a selected destination, the flow generally representing a single session using a protocol. Each packet or cell in a flow is expected to have the same routing information according to the dictates of the protocol. The specific embodiment of the MS-SAR device described here can process up to one million flows simultaneously. Each of these one million flows is processed in accordance with any one of several "application types". In each application type, the ingress MS-SAR processes the flow in accordance with one of several "ingress application types" and the egress MS-SAR processes the flow in accordance with one of several "egress application types". Figure 6 sets forth, for each application type, the related ingress application type and the related egress application type. Figure 7 sets forth the ingress and egress application types when the switch fabric is a cell-based switch fabric. Figure 8 sets forth the ingress and egress application types when the switch fabric is a packet-based switch fabric.

EXAMPLE OF APPLICATION TYPES 0, 3, 8 AND 11:

Figures 9-37 illustrate an example of how two flows of different traffic types (flow #1 is an MPLS packet, flow #2 is an ATM cell) are received onto a first line card (in this example, line card 101), simultaneously pass through the ingress MS-SAR 125 on the first line card, pass through switch fabric 105 (in this case switch fabric 105 is a cell-based switch fabric), pass onto a second line card (in this example, line card 103), simultaneously pass through an egress MS-SAR 200 on the second line card 103, and are output by the second line card 103. In this example, flow #1 and flow #2 are both communicated to first line card 101 on the same wavelength channel transmitted on the same fiber optic cable 115. Similarly, flow #1 and flow #2 are both communicated from second line card 103 on the same wavelength channel transmitted on the same fiber optic cable 118.

In Figure 9, flow #1 is an MPLS packet flow passing into an ingress line card, where the ingress line card supplies switch cells to the switch fabric. As indicated by the lower left example of Figure 7, this type of flow is identified as ingress application type 3. Ingress MS-SAR 125 therefore performs processing on flow #1 in accordance with ingress application type 3. In Figure 9, flow #2 is an ATM cell flow passing into an ingress line card, where the ingress line card supplies switch cells to the switch fabric. As indicated by the upper left example of Figure 7, this type of flow is identified as ingress application type 0. Ingress MS-SAR 125 therefore performs processing on flow #2 in accordance with ingress application type 0.

Figure 10 is a more detailed diagram of the MS-SAR devices 125 and 200 present on line cards 101 and 103. Each MS-SAR device includes an incoming SPI-4 interface block 201, a lookup engine block 202, a segmentation block 203, a memory manager block 204, a reassembly and header-adding block 205, an outgoing SPI-4 interface block 206, a per flow queue block 207, a data base block 208, a traffic shaper block 209, an output scheduler block 210, and a CPU interface block 211. In one embodiment, the blocks within dashed line 212 are realized on a first integrated circuit and the blocks within dashed line 213 are realized on a second integrated circuit. In the event the functionality of the MS-SAR is broken into two integrated circuits, a second CPU interface block 214 is provided so that CPU 106 can communicate with both integrated circuits. The MS-SAR interfaces to and uses

1 numerous other external memory integrated circuit devices 215-228 that are
2 disposed on the line card along with the MS-SAR.

3 In the example of Figure 9, the MPLS packet of flow #1 is received onto input
4 optics module 119 from fiber optic cable 115. The MPLS packet passes through
5 (see Figures 5 and 9) the input optics module 119, through serdes 121, and to multi-
6 service SONET framer 123. Framer 123 is a framer that includes a demapper for
7 separating ATM cells and MPLS packets that are received on the same wavelength
8 via input optics module 119. Framer 123 in one embodiment is a Ganges S19202
9 STS-192 POS/ATM SONET/SDH Mapper available from Applied Micro Circuits
10 Corporation, 200 Brickstone Square, Andover, MA 01810. Framer 123 outputs the
11 MPLS packet information of flow #1, along with a logical input port number indicative
12 of the physical input module from which the information was received, in 16-bit
13 pieces, onto SPI-4 bus 134. Classification engine 124 performs IP (Internet
14 Protocol) classification in the case where router 100 performs IP lookup.
15 Classification engine 124 may, in one embodiment, be a classification engine
16 available from Fast-Chip Incorporated, 950 Kifer Road, Sunnyvale, CA 94086. In
17 the particular example of Figure 9, the packet traffic (flow #1) comes into the router
18 in MPLS form and exits the router in MPLS form. Consequently, IP classification is
19 not performed in this embodiment.

20 Figure 11 is a simplified diagram of the MPLS packet 300 of flow #1 as it is
21 received onto ingress MS-SAR 125. MPLS packet 300 contains an MPLS header
22 301, and a data payload 303. An SPI-4 start delimiter 302 and an SPI-4 end
23 delimiter 304 frame the packet. MPLS header 301 includes the 20-bit MPLS label.
24 The MPLS header 301 is disposed in the packet between the layer two (L2) header
25 and the layer three (L3) IP header. The L2 header was removed by framer 123 such
26 that beginning of the packet as received onto ingress MS-SAR 125 is the MPLS
27 header.

28 Ingress MS-SAR 125 (see Fig. 10) receives the packet information in multiple
29 sixteen-byte bursts, sixteen bits at a time, via sixteen input terminals 198. When
30 incoming SPI-4 interface block 201 accumulates 64 bytes (block 201 accumulates
31 64 bytes in ingress mode and up to 80 bytes in egress mode), it sends the 64-byte
32 chunk to lookup block 202 via 64-bit wide data bus 317. In this way incoming SPI-4

1 interface block 201 breaks the packet up into 64-byte chunks and communicates the
2 chunks to lookup block 202.

3 Figure 12 illustrates the three 64-byte chunks into which this particular MPLS
4 packet 300 is broken as output from incoming SPI-4 interface block 201 to lookup
5 block 202. The incoming SPI-4 interface block 201 uses the start delimiter 302 and
6 the end delimiter 304 to identify the beginning and ending of the packet 300. The
7 incoming SPI-4 interface block 201 outputs two additional signals (an end-of-packet
8 signal and a start-of-packet signal) along with each of the 64-byte chunks. Figure 12
9 shows the end-of-packet (EOP) and start-of-packet (SOP) bits for each of the three
10 64-byte chunks.

11 Figure 13 shows information flow through ingress MS-SAR 125. CPU 106 has
12 previously placed lookup information into ingress MS-SAR 125 so that the
13 information in each MPLS packet header can be used by lookup block 202 to find: 1)
14 a particular flow ID (FID) for the flow that was specified by CPU 106, and 2) the
15 ingress application type. The ingress application type, once determined, is used by
16 other blocks of the ingress MS-SAR 125 to configure themselves in the appropriate
17 fashion to process the network information appropriately.

18 Only one traffic type (for example, MPLS packet or ATM cell) is permitted on
19 each logical input port of SPI-4 bus 135. The traffic type is assigned by CPU 106,
20 but there can only be one traffic type assigned at a given time for a given logical
21 input port. A "port table" in lookup block 202 contains, for each of the sixty-four
22 logical input ports, the allowed traffic type as specified by CPU 106. This allowed
23 traffic type information is provisioned in the port table before the flow starts by CPU
24 106.

25 Figure 14 is a conceptual diagram that shows the relationship of the information
26 in the "port table" in lookup block 202. For each logical input port there is a traffic
27 type defined. Lookup block 202 uses the incoming logical port number supplied by
28 framer 123 to lookup the traffic type allowed on the logical input port that packet 300
29 was received on. Each traffic type has a known format. Accordingly, once the traffic
30 type is known, the lookup block can locate in the header of the incoming flow the
31 particular information that will be used to lookup an associated FID found in a "FID

1 hash table". This "FID hash table" is stored in external memories 215 and 216 (see
2 Fig. 10).

3 In the present example of Figure 9, packet 300 is an MPLS packet. Lookup block
4 202 uses the traffic type found in the "port table" to locate the 20-bit MPLS label
5 within the first chunk of the packet. The MPLS label contains information similar to
6 source and destination addressing information. A hash generator in the lookup
7 block 202 uses the 20-bit MPLS label as a "hash key" to generate a "hash index".
8 This "hash index" is used as an address for the "FID hash table". The content of the
9 FID hash table location addressed contains both the FID and the ingress application
10 type. For more detailed information on how this FID lookup is performed in one
11 particular embodiment, see U.S. Patent Application Serial No. 09/779,381, filed
12 February 7, 2001, by Parruck et al. (the subject matter of which is incorporated
13 herein be reference).

14 The FID and ingress application type, once determined, are passed via 64-bit
15 wide data bus 318 (see Figure 10), one at a time, to segmentation block 203 along
16 with the consecutive 64-byte chunks of data. Segmentation block 203 stores these
17 64-bytes chunks on a per port basis. Each 64-byte chunk is stored in a buffer in
18 SRAM within segmentation block 203, and a pointer to the buffer is pushed onto a
19 FIFO of pointers for the input port. The FID and traffic type is also pushed onto the
20 FIFO with the pointer. In the present example, there are three 64-byte chunks in
21 flow #1. Segmentation block 203 calculates a cyclic redundancy check (CRC) value
22 for the data portion of the entire packet, and adds a trailer including this CRC such
23 that the trailer is at the end of the last 64-byte chunk. Segmentation block 203 adds
24 any pad that might be necessary between the end of the data of the last 64-byte
25 chunk and the trailer. Memory manager block 204 pops the FIFO of pointers such
26 that segmentation block 203 forwards the 64-byte chunks one at a time to memory
27 manager block 204 via a 128-bit wide bus 319.

28 Payload memory 217 contains a large number of 64-byte buffers, each buffer
29 being addressed by a buffer identifier (BID). Some of the 64-byte buffers of payload
30 memory 217 may be used and storing data, whereas others may not be storing data
31 (i.e., are "free"). Per flow queue block 207 maintains a linked list of the BIDs of
32 available buffers, called the "free buffer queue". The "free buffer queue" is stored in

1 external SSRAM 226 (see Figure 10). When memory manager block 204 receives
2 the first 64-byte chunk 305 of data associated with packet 300, memory manager
3 block 204 issues an “enqueue” command via enqueue command line 320 to per flow
4 queue block 207. This constitutes a request for the BID of a free buffer. Per flow
5 queue block 207 pops the free buffer queue to obtain the BID of a free buffer, and
6 forwards that BID to memory manager block 204 via lines 321. Memory manager
7 block 204 then stores the 64-byte chunk 305 (see Fig. 12) of data for packet 300 in
8 the buffer in payload memory 217 identified by the BID. The writing of the 64-byte
9 chunk of data is indicated in Figure 13 by the upward pointing heavy arrow 322 that
10 extends toward payload memory block 217. For additional details on the operation
11 of the memory manager block 204 and how it stores information into payload
12 memory 217, see U.S. Patent Application Serial No. 09/779,381, filed February 7,
13 2001, by Parruck et al. (the subject matter of which is incorporated herein by
14 reference).

15 Per flow queue block 207 also maintains a linked list (i.e., a “queue”) of the BIDs
16 for the various 64-byte chunks of each flow that are stored in payload memory 217.
17 Such a linked list is called a “per flow queue”. Figure 15A illustrates how per flow
18 queue block 207 builds the linked list for flow #1. Each linked list has a head pointer
19 and a tail pointer. A list of the head pointers for the one million possible FIDs is
20 stored in external SRAM 228 (see Figure 10). A list of the tail pointers for the one
21 million possible FIDs is stored in external SRAM 227.

22 The head pointer for the FID of flow #1 is set to point to the memory location that
23 stores the BID where the first chunk 305 is stored in payload memory 217. Because
24 there is only one chunk (C1) in the linked list for flow #1, the tail pointer is set to
25 point to the same location. This first chunk 305 of flow #1 is the start of packet 300
26 as identified by an SOP bit.

27 Next, the second 64-byte chunk 306 is received by the memory manager block
28 204. The memory manager block 204 again issues an enqueue command to per
29 flow queue block 207, again obtains a BID of a free buffer, and then writes the
30 second chunk 306 into the payload memory at the location identified by the BID.
31 The per flow queue block 207 pops the free buffer queue to remove the now-used
32 BID from the free buffer queue, and adds the BID for the second chunk 306 to the

1 linked list for the FID of flow #1. As illustrated in Figure 15B, the BID of the second
2 chunk 306 is added to the linked list for the FID of flow #1.

3 Next, the third 64-byte chunk 307 is received by the memory manager block 204.
4 The same sequence of storing the chunks into payload memory and adding the BID
5 to the linked list of flow #1 is carried out. Figure 15C illustrates the per flow queue
6 for flow #1 including pointers to these three chunks. In the case of chunk 307, this
7 chunk is the last chunk of the packet as indicated by the EOP bit. The linked list for
8 flow#1 is therefore complete.

9 Figure 16 is a diagram of a memory location in FID dequeue memory 228 (see
10 Figure 10) that stores the head pointer for a flow. Per flow queue block 207 stores in
11 the location, along with the BID head, other information it receives from memory
12 manager block 204 relating to the chunk. This information includes an EOP bit that
13 indicates whether the chunk is the last chunk of a packet, an SOP bit that indicates
14 whether the chunk is the first chunk of a packet, and the ingress application type.
15 Per flow queue block 207 also stores with the head pointer EFCI, CLP and OAM
16 bits. These bits were extracted by lookup block 202 from the header. There is one
17 memory location such as the one illustrated in Figure 16 for each per flow queue. In
18 the example of Figure 15C, chunk C1 is pointed to by a head pointer in one such
19 memory location.

20 Figures 17 and 18 are diagrams of two memory locations in FID enqueue
21 memory 227 (see Figure 10) that store a tail pointer. Per flow queue block 207
22 stores in these locations, along with a pointer to the BID tail, other information on the
23 chunk including the output port number from which the data will eventually be output,
24 the size of the per flow queue, and an indication of a time to live (a TTL bit). There
25 is one such pair of memory locations for each per flow queue. In the example of
26 Figure 15C, chunk C3 is pointed to by information in a pair of two such memory
27 locations.

28 Figure 19 is a diagram of a memory location in SRAM 226 that stores the pointer
29 for a chunk in a queue between the head and tail. Per flow queue block 207 stores
30 in this location a pointer to the next BID in the per flow queue. The location also
31 stores an EOP indication and an SOP indication. In the example of Figure 15C,
32 there is one such memory location for chunk C2.

In the example of Figure 9, a second flow (flow #2) of a different traffic type (ATM cell) is received onto the same line card 101 that flow #1 was. In this example, flow #2 is also received via the same fiber optic cable 115 and is modulated onto the same wavelength channel that flow #1 was. The 53-byte ATM cell of flow #2 passes through the same optics module 119, serdes 121 and framer 123. Framer 123 removes the fifth byte (the "HEC" Header Error Control byte) of the ATM header and places the remaining 52-byte ATM cell onto the 16-bit SPI-4 bus 134. This information passes through classification engine 124 to SPI-4 bus 135 such that the 52-byte ATM cell is received onto the incoming SPI-4 interface block 201 of ingress MS-SAR 125 via sixteen SPI-4 input terminals 198 (see Fig. 10). Incoming SPI-4 interface block 201 receives the 52-byte ATM cell in multiples of 16-byte bursts, 16 bits at a time.

Figure 20 is a simplified diagram of the ATM cell 308 as received onto ingress MS-SAR 125. ATM cell 308 is 52-bytes long. ATM cell 308 includes an ATM header 309 and a data payload 311. An SPI-4 bus start delimiter 310 and end delimiter 312 frame the ATM cell. Incoming SPI-4 interface block 201 receives the ATM cell and supplies it to lookup block 202, sixty-four bits at a time, via 64-bit wide data bus 317.

The internal data path within the MS-SAR from incoming SPI-4 interface block 201, through lookup block 202, and to segmentation block 203 is 64-bits (8-bytes) wide. Incoming SPI-4 interface block 201 therefore adds a 4-byte pad 314 (see Figure 21) to pad the 52-byte ATM cell up to 56-bytes (a multiple of 8 bytes) before sending the ATM cell to lookup block 202.

Figure 21 is a diagram of ATM cell 308 as it is output from incoming SPI-4 interface block 201 to lookup block 202. Because the entire ATM cell 308 is contained in one 64-byte chunk, the SOP and EOP signals output by the incoming SPI-4 interface block 201 indicate both the start of packet (in this case the packet is a cell) and end of packet (in this case the packet is a cell).

Lookup block 202 receives the 56-byte chunk for ATM cell 308, and from the logical input port number looks up the traffic type from the "port table" (see Figure 14). The "port table" indicates that the traffic type is ATM cells. Lookup block 202 uses the traffic type to locate the 12-bit VPI and 16-bit VCI fields in the ATM cell

1 header. The hash generator of lookup block 202 uses the located VPI and VCI
2 information as a “hash key” to create a “hash index”. This “hash index” is then used
3 to lookup the FID and ingress application type in the “FID table” stored in external
4 memories 215 and 216.

5 Once the FID and the ingress application type for flow #2 are determined, these
6 values are passed from lookup block 202 to segmentation block 203 via 64-bit bus
7 318 along with the 56-byte chunk. Segmentation block 203 adds an additional 8-
8 byte pad to pad the ATM cell up to 64-bytes, stores the 64-byte chunk into an SRAM
9 (not shown) in segmentation block 203, and pushes a pointer to that chunk onto its
10 FIFO of pointers for that input port. Memory manager block 204 pops the FIFO of
11 pointers such that segmentation block 203 supplies the 64-byte chunk from its
12 SRAM to memory manager block 203.

13 Figure 22 is a simplified diagram of the 64-byte chunk as output from
14 segmentation block 203 onto 128-bit bus 319. This 64-byte cell contains the 4-byte
15 pad 314 added by the incoming SPI-4 interface block 201 and the additional 8-byte
16 pad 315 added by segmentation block 203. When memory manager block 204
17 receives the 64-byte chunk, memory manager block 204 issues an enqueue
18 command via command line 320 to per flow queue block 207, obtains in return a
19 BID, and then stores the 64-byte chunk in payload memory 217 at the location
20 identified by the BID. Per flow queue block 207 pops the free buffer queue, thereby
21 removing the now-used BID from the free buffer queue, and adds the BID to a linked
22 list for flow #2.

23 Figure 23 is a diagram that illustrates the two linked lists. The linked list for flow
24 #1 (FID1) has BIDs for three chunks C1, C2 and C3, whereas the flow for flow #2
25 (FID2) has BIDs for one chunk C1. The chunk for the ATM cell is indicated as being
26 both the “start of packet” as well as the “end of packet”.

27 Once the linked lists (queues) for flow #1 and flow #2 are formed, the linked lists
28 are popped (i.e., dequeued) in a particular way and order such that their associated
29 chunks in payload memory 217 are output from ingress MS-SAR 125 in a desired
30 fashion, taking into account user-programmable policing/metering parameters.

31 User-programmable parameters may include, for example, burst size, peak cell rate,
32 and sustainable cell rate.

1 The dequeue process starts with data base block 208 determining an output port
2 to service using a port calendar. This port calendar is located within data base block
3 208. Once the output port to service is selected, data base block 208 supplies the
4 output port to traffic shaper block 209 and to output scheduler block 210. The traffic
5 shaper block 209 and the output scheduler block 210 supply flow IDs back to data
6 base block 208. Data base block 208 selects one of the two flow IDs to dequeue for
7 that port. Data base block 208 gives priority to traffic shaper output over output
8 scheduler output such that only if the shaper has no FID to output for a given port
9 will the output scheduler be allowed to schedule an FID for that port. Either traffic
10 shaping is performed on an individual flow by traffic shaper block 209, or output
11 scheduling is performed on the flow by output scheduler block 210.

12 Traffic shaper block 209 performs traffic shaping on a per flow basis for up to one
13 million flows. On a particular flow, either a single leaky bucket shaping scheme is
14 used, a dual leaky bucket shaping scheme is used, or no shaping at all is used. The
15 shaper selects a peak rate or a sustained rate per flow ID depending on an
16 accumulated credit. Up to 1024 different programmable rates can be shaped
17 simultaneously. In one mode, the rate of incoming traffic with a particular flow ID is
18 measured (i.e., metered) and the flow ID of a violator is marked if the measured rate
19 is above a specific threshold programmed by CPU 106 for the flow ID. Up to 1024
20 different programmable thresholds can be metered simultaneously.

21 Output scheduler block 210 uses a weighted round-robin scheme to select a
22 quality of service of the selected port. Once the quality of service is selected, a flow
23 ID is selected based on a round-robin scheme.

24 When data base block 208 receives a flow ID from either traffic shaper block 209
25 or output scheduler block 210, data base block 208 generates a request to per flow
26 queue block 207 to issue a dequeue command to dequeue the flow ID. Per flow
27 queue block 207 accesses the per flow queue of the flow ID, determines the next
28 BID to dequeue, and outputs the BID in the form of a dequeue command to memory
29 manager block 204.

30 Per flow queue block 207 is programmable to cause the linked list of one flow to
31 be output multiple times (i.e., multicast), each time with a different flow ID.

32 Multicasting is performed by replicating the linked list of a flow, each replica linked

1 list having its own flow ID. Per flow queue block 207 is also programmable to cause
2 multiple flows to be combined into one flow (i.e., tunneling). Tunneling is performed
3 by linking the tail of one linked list to the head of another so as to create one large
4 composite linked list, the composite linked list having one flow ID.

5 For additional structural and operational details of one particular embodiment of
6 traffic shaper block 209, one particular embodiment of output scheduler block 210,
7 one particular embodiment of data base block 208, and one particular embodiment
8 of per flow queue block 207, see U.S. Patent Application Serial No. 09/779,381, filed
9 February 7, 2001, by Parruck et al. (the subject matter of both of which is
10 incorporated herein by reference).

11 In the present example, the linked list of flow #1 (see Figure 22) is dequeued first.
12 Per flow queue block 207 pops the BID of the first chunk C1 off the FID1 linked list
13 for flow #1 and forwards that BID to memory manager block 204 in a dequeue
14 command. The dequeue command contains the BID of the 64-byte chunk to be
15 read out of payload memory 217, as well as the ingress application type, an EOP bit,
16 an SOP bit, and the output port number (one of 64 logical output ports on SPI-4 bus
17 107). The dequeue command is sent via dequeue command line 323. The BID is
18 sent via BID lines 321. Per flow queue block 207 adds the now available BID to the
19 free buffer queue in external memory 226.

20 In response to receiving the dequeue command, memory manager block 204
21 retrieves the first chunk C1 identified by the BID and outputs that first chunk C1 to
22 reassembly block 205 via 128-bit data bus 324. Memory manager block 205
23 supplies to reassembly block 205 other control information including the FID of
24 chunk C1, the SOP and EOP bits, the ingress application type being performed on
25 flow #1, and a logical output port ID (PID) identifying the one of the 64 logical output
26 ports on SPI-4 bus 107 to which the chunk will be sent.

27 Reassembly block 205 uses the ingress application type to determine what type
28 of action it should perform on the associated chunk. In the present example, flow #1
29 is processed in accordance with ingress application type 3. External memory 218
30 (see Fig. 10) contains a "header table" 327. For each FID, the CPU 106 has stored
31 beforehand in header table 327 a switch header that reassembly block 205 is to
32 append to the data of the chunk before sending the data on to the switch fabric.

1 Accordingly, reassembly block 205 uses the FID to lookup in "header table" 327 the
2 "switch header" placed there for this FID by CPU 106. The switch header can be
3 either 8-bytes or 16-bytes depending on the requirements of the switch fabric. As
4 explained above, the "switch header" contains the "virtual port number" of the
5 particular egress line card that the switch fabric 105 wants the switch cell routed to.
6 In the present example of Figure 9, the egress line card is line card 103.
7 Accordingly, CPU 106 has placed the "virtual port number" of line card 103 into the
8 "switch header" associated with flow #1.

9 Reassembly block 205 also provides a special 4-byte "Azanda header". The
10 Azanda header is embedded in the switch cell as the last four bytes of the switch
11 header. This Azanda header contains information on how the switch cell should be
12 processed by the particular egress MS-SAR that receives the switch cell (in this
13 case MS-SAR 200). The Azanda header includes the egress application type to be
14 used to process the switch cell in the egress MS-SAR. The Azanda header also
15 includes the FID of the cell, an SOP bit, an EOP bit, and quality of service
16 information.

17 When reassembly block 205 receives the 64-byte chunk from memory manager
18 block 204, it stores the 64-byte chunk into a dual-port DATA_SRAM and stores the
19 "switch header" (including the "Azanda switch header") into a HDR_SRAM. The
20 DATA_SRAM (not shown) and the HDR_SRAM (now shown) are part of the
21 reassembly block 205. A pointer that points to the data in DATA_SRAM and to the
22 header in HDR_SRAM is pushed onto a queue for the particular logical output port.
23 In the present example, where the MS-SAR is operating in the ingress mode, there
24 is only one logical output port (i.e., the output of line card 101), consequently there is
25 only one output port queue. The pointer is therefore placed on this queue. The
26 output port queue is maintained in a Q_FIFO (not shown) that is also part of the
27 reassembly block 205.

28 Figure 24 is a conceptual diagram of a "port calendar" located in reassembly
29 block 205. Figure 25 is a diagram of a "port empty" register located in reassembly
30 block 205. Figure 26 is a diagram of a "port full" register located in reassembly block
31 205. Reassembly block 205 uses the port calendar and registers of Figures 24-26 to
32 determine how to pop the output port queues in the Q_FIFO and move the

1 associated chunks to an OUTPUT_DATA_FIFO. The OUTPUT_DATA_FIFO (not
2 shown) is located within reassembly block 205 and serves as the interface through
3 which data flows to outgoing SPI-4 interface block 206.

4 Outgoing SPI-4 interface block 206 periodically sends port full status information
5 that is used to load the “port full” register to keep reassembly block 205 updated on
6 how busy the various output logical ports are. Memory manager block 204
7 periodically sends empty status information that is used to load the “port empty”
8 register to keep reassembly block 205 updated on whether there is any more
9 information coming from the memory manager 205 for each of the various logical
10 output logical ports. In each of the “port empty” and “port full” registers, there is one
11 bit for each of the respective sixty-four logical output ports of ingress MS-SAR 125.

12 Reassembly block 205 steps through the rows of the port calendar of Figure 24,
13 one by one. For each row, reassembly block 205 reads the port ID (PID) from the
14 last field of the row and examines the bit corresponding to that port ID in each of the
15 “port empty” and “port full” registers. If the bits in these two registers indicate that
16 the port is neither full (i.e., that port has all the data it can handle) nor empty (no data
17 is available for that port), then reassembly block 205 pops the queue in Q_FIFO for
18 that port ID. The popping of the queue causes reassembly block 205 to retrieve one
19 64-byte data chunk from DATA_SRAM and its associated header from HDR_SRAM.
20 (In ingress mode, all chunks are 64 bytes but in egress mode EOP chunks can be
21 less than 64 bytes). The header is then appended to the data chunk as the complete
22 “switch cell” is pushed onto the OUTPUT_DATA_FIFO. Outgoing SPI-4 interface
23 block 206 pops the OUTPUT_DATA_FIFO, thereby causing the “switch cell” to be
24 supplied to outgoing SPI-4 interface block 206 in 64-bit pieces via 64-bit bus 325
25 (see Figure 10). The logical output port is provided to outgoing SPI-4 interface block
26 206 along with the switch cell. Outgoing SPI-4 interface block 206 in turn outputs
27 the switch cell onto SPI-4 bus 107 in 16-bit pieces via sixteen output terminals 326
28 of ingress MS-SAR 125. The SPI-4 bus 107 includes an SOP delimiter and an EOP
29 delimiter between data bursts on the SPI-4 bus. Data bursts on the SPI-4 bus are
30 transmitted in multiples of 16 bytes.

31 In this example, MS-SAR 125 is in the ingress mode, so there is only one logical
32 output port. Consequently there is only one entry in the port calendar of Figure 24.

1 When a row entry is processed having a jump bit of “1”, processing of rows returns
2 to the first row. In the example of Figure 24 where there is only one logical output
3 port, there is a “1” in the jump column of the first row entry for port ID number one.
4 Accordingly, the one and only output queue in the Q_FIFO is popped. The circular
5 arrow in Figure 24 illustrates this return in the present example where there is only
6 one logical output port (PID 1).

7 Figure 27 illustrates information flow for flow #1 out of ingress MS-SAR 125. The
8 heavy upwardly pointing arrow represents the dequeuing operation being performed
9 by per flow queue block 207 including the passing of BIDs and associated
10 parameters to memory manager block 204. The heavy arrow extending from
11 payload memory block 217 to the right represents the flow of 64-byte chunks
12 through reassembly block 205 and to the switch fabric.

13 Per flow queue 207 continues issuing dequeue commands in the desired order
14 such that the associated various 64-byte chunks of the two flows (see Figure 23) are
15 read out of payload memory 217 and such that corresponding switch cells are
16 generated. The switch cells are output to switch fabric 105 from outgoing SPI-4
17 interface block 206 via SPI-4 bus 107.

18 Figure 28 illustrates the three switch cells that correspond to 64-byte chunks C1,
19 C2 and C3 of packet 300. The switch headers of each of the three cells contain the
20 “virtual port number” of line card 103 to which the cells are to be routed. The
21 Azanda header of each switch cell contains the egress application type (egress
22 application type 11) that the egress MS-SAR 200 is to perform on the switch cells of
23 flow #1.

24 Figure 29 illustrates the switch cell that corresponds to the one and only 64-byte
25 chunk for ATM cell 308. The switch header of this switch cell contains the “virtual
26 port number” of line card 103 to which the cell is to be routed. The Azanda header
27 of the switch cell contains the egress application type (egress application type 8) that
28 the egress MS-SAR 200 is to perform on the switch cell of flow #2.

29 Switch fabric 105 receives the switch cells from SPI-4 bus 107 and uses the
30 “virtual port number” embedded in the switch headers to route the switch cells to line
31 card 103. The switch cells are received on the SPI-4 bus 110 (see Fig. 9) of line
32 card 103. SPI-4 bus 110 is coupled to the SPI-4 bus input terminals (see input

terminals 198 of Figure 10) of egress MS-SAR device 200 of line card 103.

Operation of egress MS-SAR 200 is explained below in connection with Figure 10.

Figure 30 illustrates the general flow of information into egress MS-SAR 200.

The first incoming switch cell of packet 300 is read into incoming SPI-4 interface block 201 and is passed to lookup block 202. Lookup block 202 need not, however, perform a lookup using hashing as was done in the ingress mode described above. Rather, because there is an Azanda header embedded in the switch cells, and because that Azanda header contains the FID and the egress application type to be performed by the egress MS-SAR, lookup block 202 merely locates this Azanda header and extracts the embedded information.

Figure 31 illustrates the first switch cell for packet 300 of flow #1. The Azanda header is the last four bytes of the switch header. Although the size of the switch header can differ from switch fabric to switch fabric, all switch cells for a given switch fabric have the same size. CPU 106 writes the size of the switch header for its switch fabric into a "number of bytes of switch header" field in each row of the "port table" (see Figure 14). This number is the same for all rows of the "port table". Lookup block 202 in egress MS-SAR 200 reads the "number of bytes of switch header" field in the row of the port table in the egress MS-SAR 200 for the logical input port from which the switch cells were received. From this number of bytes of the switch cell, lookup engine block 202 locates the last four bytes which are the Azanda header. From the Azanda header, lookup block 202 retrieves the FID and the egress application type to be performed by egress MS-SAR 200. Lookup block 202 removes the switch header and passes the remaining 64-bytes to segmentation block 203 along with the FID and egress application type. The egress application type is supplied along with the 64-byte chunk so that the other blocks that receive the 64-byte chunk will perform the correct type of processing on the associated 64-bytes. The FID in the example of Figure 9 is flow #1. The egress application type in the example of Figure 9 for flow #1 is egress application type 11.

In accordance with egress application type 11, segmentation block 203 does no segmenting per se but rather forwards the 64-byte chunk to memory manager block 204 via the per-port FIFO mechanism described above. Memory manager block 204 issues an enqueue command via enqueue command line 320, and stores the 64-

1 byte chunk in payload memory 217. Per flow queue block 207 adds the BID to the
 2 per flow queue for flow #1. The same process flow occurs for the second and third
 3 switch cells associated with flow #1.

4 Figure 30 represents the storing of 64-byte chunks into payload memory 217 by
 5 the heavy upward pointing arrow that extends toward payload memory block 217.
 6 The heavy downward pointing arrow that extends toward per flow queue block 207
 7 represents the enqueue command and the building of the three BID linked list for
 8 flow #1.

9 The switch cell for the ATM cell (see Fig. 29) is received onto egress MS-SAR
 10 200 (see Fig. 9) in similar fashion. As in the case for the switch cells for the packet,
 11 lookup block 202 uses the logical input port number supplied by incoming SPI-4
 12 interface block 201 to access the appropriate row in the "port table" (see Fig. 14) in
 13 the lookup block 202. From the "number of bytes of switch header" field in that row,
 14 lookup block 202 locates the "Azanda header" in the incoming switch cell. The
 15 Azanda header contains the egress application type and the FID. Lookup block 202
 16 removes the switch header, and forwards the remaining 64-byte chunk to
 17 segmentation block 203 along with the FID and egress application type to be
 18 performed by egress MS-SAR 200. The 64-byte chunk in this case contains the 4-
 19 byte ATM header, the 48-byte ATM payload, and twelve bytes of pad.

20 In the example of Figure 9, the FID is flow #2. In the example of Figure 9, the
 21 egress application type for flow #2 is egress application type 8. In accordance with
 22 egress application type 8, segmentation block 203 does not perform segmentation
 23 per se but rather passes the 64-byte chunk and the egress application type to
 24 memory manager block 204 via the per-port FIFO mechanism described above.
 25 Memory manager block 204 stores the 64-byte chunk. Per flow queue block 207
 26 adds the BID for this chunk and parameters associated with the chunk to a per flow
 27 queue for flow #2.

28 Once the 64-byte chunks from the various switch cells are stored in payload
 29 memory 217 and the linked lists for flow#1 and flow #2 are built, then the linked lists
 30 can be popped (i.e., "dequeued") and the 64-byte chunks output in various orders.
 31 For example, output scheduler block 210 can schedule the popping of the queues to
 32 control flow to output ports on a flow ID basis.

Figure 32 illustrates the general flow of information out of egress MS-SAR 200 in one scenario when the linked lists for flow #1 and flow #2 are dequeued. As each BID is dequeued from a per flow queue, the associated 64-byte chunk is read from payload memory 217 and is supplied via 128-bite wide data bus 324 to reassembly block 205. As explained above in connection with the ingress mode, reassembly block 205 maintains one reassembly queue for each of the 64 logical output ports. If two flows share the same logical output port, then the entire linked list for one flow must be dequeued before the linked list for the next flow having the same logical output port is dequeued. In the example of Figure 8, flow #1 and flow #2 have different logical output ports. The dequeue command as received by memory manager block 204 from per flow queue block 207 contains the port ID (PID) of the one of the 64 logical output ports of line card 103.

In the present example, the three-chunk linked list for flow #1 is dequeued first, one chunk at a time. Each 64-byte chunk is read into the DATA_SRAM of reassembly block 205 and a pointer to that chunk is stored in the Q_FIFO in the queue for the logical output port of flow #1. The first chunk contains the MPLS packet header. In an MPLS manipulation operation, reassembly block 205 can add (i.e., "push") MPLS labels to, or remove (i.e., "pop") MPLS labels from, the MPLS header found in the first chunk.

Consider the situation in Figure 9 where a first MPLS label for communication across the network coupled to router 100 via fiber optic cable 115 is to be replaced with a second MPLS label for communication across the network coupled to router 100 via fiber optic cable 118. Reassembly block 205 determines the type of MPLS manipulation to perform using the FID to lookup MPLS information that CPU 106 stored previously in the external "header table" 327 (see Figure 10) of the egress MS-SAR 200. For each FID, header table 327 contains MPLS manipulation information.

Figure 33 is a more detailed diagram of an entry for a FID in "header table" 327 of Figure 10. The entry for each FID includes two 72-bit words. The first word contains an eight-bit control word CONTROL_WORD_1, a first 32-bit field that can contain either a first MPLS label or an ATM VPI/VCI header, and a second 32-bit field that can contain a second MPLS label. The second word contains a second eight-bit

control word CONTROL_WORD_2, a first 32-bit field that can contain a switch header, and a second 32-bit field that holds the 32-bit Azanda header. The first two bits of CONTROL_WORD_1 indicate the header type: "00" indicates MPLS type, "10" indicates ATM type, and "01" indicates L2 Ethernet type. The next three bits of CONTROL_WORD_1 indicate the type of MPLS manipulation to perform: "000" indicates replace MPLS label and decrement the 8-bit MPLS time to live (TTL) field, "001" indicates push one MPLS label and decrement the original MPLS TTL, "010" indicates push two MPLS labels and decrement the original MPLS TTL, "011" indicates pop one MPLS label and decrement the original MPLS TTL, and "100" indicates pop two MPLS labels and decrement the original MPLS TTL. The next two bits of CONTROL_WORD_1 indicate the number of valid header words for that FID in the header table 327. The last bit of CONTROL_WORD_1 indicates that L2 stripping and CRC checking is enabled. The first five bits of CONTROL_WORD_2 indicate the number of L2 header bytes to remove if type is L2 Ethernet.

In the example of Figure 9, the entry in header table 327 for flow #1 indicates that the header type is MPLS and that information in the MPLS field (see Figure 33) determines the action. In the present example, the MPLS label is to be replaced. Reassembly block 205 therefore replaces the MPLS label in the first 64-byte chunk for flow #1 with the MPLS label found in an MPLS label field of the first word of the entry for flow #1 in header table 327. Reassembly block 205 also decrements the original TTL number.

As described above in connection with Figures 24-26, reassembly engine 205 cycles through the entries in its port calendar (see Fig. 24). The port calendar is provisioned beforehand by CPU 106 so that there is one entry for the port of flow #1 and another entry for the port of flow #2. If when reassembly block 205 services one of the rows of the port calendar the full and empty registers indicate that outgoing SPI-4 interface block 206 is ready for data for that output port, then the queue for that output port is popped such that one 64-byte chunk stored in DATA_SRAM is combined with a header from HDR_SRAM if appropriate and the result is pushed onto the OUTPUT_DATA_FIFO within reassembly block 205.

As the various 64-byte chunks of the packet are received from memory manager block 204, reassembly block 205 uses a CRC engine to calculate a CRC value on

1 the data. The CRC value calculated on only part of the data is called a “partial
2 CRC”. The partial CRC value is changed as each additional piece of data passes
3 through the reassembly block 205. When the last of the data (as indicated by the
4 EOP bit) passes through the CRC engine, the CRC value is a “complete CRC”. The
5 complete 32-bit CRC is compared with the 32-bit CRC value found in the AAL5
6 trailer of the last 64-byte chunk of the packet. If the two CRC values match, then the
7 data was not corrupted.

8 The AAL5 trailer also contains sixteen-bit “length of payload” field that indicates
9 the length of the payload of the packet being reassembled. Reassembly block 205
10 maintains a partial packet byte count value for the packet. As each 64-byte chunk of
11 the packet is received from memory manager block 204, reassembly block 205 adds
12 the number of bytes received to the partial packet byte count value. After the last
13 64-byte chunk for the packet has been received and the partial packet byte count
14 value has been added to, reassembly block 205 compares the now complete packet
15 byte count value with the “length of payload” value from the AAL5 trailer. The two
16 should match.

17 After checking the CRC and the byte count, reassembly block 205 removes the
18 trailer and any padding from the last 64-byte chunk, and pushes the last chunk onto
19 the OUTPUT_DATA_FIFO.

20 Outgoing SPI-4 interface block 206 receives the pieces of the packet, piece by
21 piece (each piece being 64-bytes), from OUTPUT_DATA_FIFO and supplies those
22 pieces of the packet of flow #1 to framer 123. As framer 123 is ready to receive
23 another piece of the packet being sent out on a particular output port, outgoing SPI-4
24 interface block 206 sends it one via 16-bit wide SPI-4 bus 141. Outgoing SPI-4
25 interface block 206 supplies bursts and start and end delimiters in accordance with
26 SPI-4 protocol to demarcate the packet. Framer 142 receives the packet from SPI-4
27 bus 141 and uses the start delimiters and the end delimiters to frame the packet.
28 The packet passes through transmit serdes 143 and optics module 144 and onto
29 fiber optic cable 118. In this way, 64-byte chunks of the packet of flow #1 go through
30 reassembly block 205 one by one, the header of the first chunk being modified in
31 accordance with egress application type 11, the trailer on the last chunk being
32 checked and then removed. Accordingly, the term “reassembly” here does not

1 necessarily require that the entire packet exist in reassembled form in the
2 reassembly block or even on the MS-SAR. Rather the reassembly block outputs the
3 various chunks of the packet sequentially such that the packet is transmitted in
4 reassembled form on the fiber optic cable.

5 Figure 34 illustrates the three 64-byte chunks of flow #1 as passed through
6 reassembly block 205 to outgoing SPI-4 interface block 206. The chunks are
7 illustrated as 64-byte chunks rather than 64-bit pieces for ease of illustration.

8 Figure 35 illustrates the MPLS packet of flow #1 as output from SPI-4 interface
9 block 206 to framer 142.

10 The 64-byte chunk for the ATM cell of flow #2 in this example is dequeued from
11 payload memory 217 in similar fashion. Per flow queue block 207 issues a dequeue
12 command to memory manager block 204 that indicates the BID of the chunk, and
13 the logical output port from which the flow is to be transmitted. Reassembly block
14 205 uses the FID of flow #2 to access the FID entry in "header table" 327 (see
15 Figure 9). The header type bits of CONTROL_WORD_1 of the FID entry for flow #2
16 indicates that the header type is ATM. Reassembly block 205 therefore replaces
17 (i.e., "translates") the VPI/VCI field in the ATM header of the chunk in HDR_SRAM
18 with the VPI/VCI information stored in the FID entry for flow #2. Reassembly block
19 205 may also be programmed by information in header table 327 to change the CLP
20 (Cell Loss Priority bit in the ATM header that indicates a high or low priority for the
21 ATM cell).

22 Reassembly block 205 receives the 64-byte chunk, stores it into DATA_SRAM,
23 and pushes a pointer to the 64-bytes onto a queue for the logical output port for flow
24 #2. Reassembly block 205 services the ports in its port calendar as explained
25 above, combines the data portion in the DATA_SRAM with the translated ATM
26 header portion in HDR_SRAM, and loads the completed ATM cell of flow #2 into the
27 OUTPUT_DATA_FIFO in reassembly block 205. The last eight bytes of the 64
28 bytes in DATA_SRAM are pad. Reassembly block 205 removes this 8-byte pad by
29 not reading it into the OUTPUT_DATA_FIFO. Outgoing SPI-4 interface block 205
30 pops the OUTPUT_DATA_FIFO such that reassembly block 205 supplies the
31 resulting 56-byte chunk to outgoing SPI-4 interface block 205. Figure 36 illustrates

1 the 56-byte chunk (ATM cell plus four bytes of pad) as output from reassembly block
2 205.

3 Outgoing SPI-4 interface block 206 removes the last four bytes of pad when
4 sending the ATM cell out. It does this using byte enables. Outgoing SPI-4 interface
5 block 206 outputs start and end delimiters to demarcate the 52-byte ATM cell in
6 accordance with the SPI-4 bus protocol. The entire 52-byte cell is output in one
7 burst. Figure 37 illustrates the ATM cell of flow #2 as output from outgoing SPI-4
8 interface block 206 of egress MS-SAR 200 via 16-bit wide SPI-4 bus 141.

9 Framer 142 receives the ATM cell from SPI-4 bus 141 and uses the start and end
10 delimiters to frame the ATM cell. Framer 142 then supplies the ATM cell to transmit
11 serdes 143. Serdes 143 in turn supplies the ATM cell in analog serial fashion to
12 optics module 144. Optics module 144 outputs the ATM cell onto fiber optic cable
13 118. In the example of Figure 9, there is only one physical output port (i.e., the fiber
14 optic cable 118). Framer 142 therefore maps the logical output ports for flow#1 and
15 flow #2 on SPI-4 bus 141 to the same serdes 143 for transmission on the same
16 wavelength channel on the same fiber optic cable 118.

18 EXAMPLE OF APPLICATION TYPES 5, 6, 14 AND 13:

19 Figure 38 illustrates an example of how two flows of different types (flow #1 is a
20 flow of three AAL5 cells, flow #2 is an MPLS packet) are received onto a first line
21 card (in this example, line card 101), simultaneously pass through ingress MS-SAR
22 125 on the first line card, pass through switch fabric 105 (in this example switch
23 fabric 105 is a packet-based switch fabric), pass onto a second line card (in this
24 example, line card 103), simultaneously pass through egress MS-SAR 200 on the
25 second line card 103, and are output by second line card 103. Flow #1 exits line
26 card 103 as an MPLS packet, whereas flow #2 exists line card 103 as four AAL5
27 cells. In this example, flow #1 and flow #2 are both communicated to first line card
28 101 on the same wavelength channel transmitted on the same fiber optic cable 115.
29 Similarly, flow #1 and flow #2 are both communicated from second line card 103 on
30 the same wavelength channel and on the same fiber optic cable 118.

31 In Figure 38, flow #1 is a flow of AAL5 cells flow passing into an ingress line card,
32 where the ingress line card supplies a packet to a packet-based switch fabric. As

1 indicated by the middle left example of Figure 8, this type of flow is identified as
 2 ingress application type 5. The various blocks within ingress MS-SAR 125 therefore
 3 perform processing on flow #1 in accordance with ingress application type 5. In
 4 Figure 38, flow #2 involves an MPLS packet flow passing into an ingress line card,
 5 where the ingress line card supplies a packet to a packet-based switch fabric. As
 6 indicated by the bottom left example of Figure 8, this type of flow is identified as
 7 ingress application type 6. The various blocks within ingress MS-SAR 125 therefore
 8 perform processing on flow #2 in accordance with ingress application type 6.

9 Ingress application type 5 processing on flow #1 is explained in connection with
 10 Figures 38, 5, 39 and 10. The first AAL5 cell of flow #1 is received via fiber optic
 11 cable 115 (see Figure 5) and passes through optics module 119, serdes 121, framer
 12 123 and classification engine 124. The 52-byte AAL5 cell is received onto ingress
 13 MS-SAR 125 in 16-bit pieces via SPI-4 bus 135 via terminals 198 (see Fig. 10). The
 14 AAL5 cell includes a 4-byte ATM header and a 48-byte payload. Incoming SPI-4
 15 interface block 201 appends a 4-byte pad to the end of the 52-byte AAL5 cell before
 16 forwarding the resulting 56-byte chunk 400 (see Fig. 39) to lookup block 202.
 17 Figure 39 illustrates the 56-byte chunk 400 as it is output from incoming SPI-4
 18 interface block 201.

19 As explained above in connection with the example of Figure 9, lookup block 202
 20 uses the input port type supplied via bus 317 to lookup the associated traffic type in
 21 the "port table" of lookup block 202. In this case, the traffic type is "ATM cells".
 22 From the traffic type, lookup block 202 locates the VPI and VCI fields in the ATM cell
 23 header. Lookup block 202 uses the VPI and VCI information to lookup in the "FID
 24 table" (in external memories 215 and 216) the FID and the ingress application type.
 25 Lookup block 202 forwards the FID and ingress application type to the other blocks
 26 of the MS-SAR. Lookup block 202 removes the 4-byte ATM header and the 4-byte
 27 pad, and forwards the remaining 48-byte chunk 401 of data to segmentation block
 28 203 via bus 318. Figure 39 illustrates the 48-byte chunk 401 as it is output from
 29 lookup block 202.

30 Segmentation block 203 adds a 16-byte pad to the 48-byte chunk to form a 64-
 31 byte chunk 402. Memory manager block 204 reads the 64-byte chunk 402 from
 32 segmentation block 203 and stores 64-byte chunk 402 into a buffer in payload

1 memory 217 using the enqueue command mechanism set forth above in connection
 2 with the example of Figure 9. The per flow queue for flow #1 at this point involves
 3 just one buffer. The other two AAL5 cells of flow #1 are processed in similar fashion
 4 such that corresponding 48-byte chunks 403 and 404 are queued in the per flow
 5 queue for flow #1.

6 Ingress application type 6 processing on flow #2 is explained in connection
 7 with Figures 38, 5, 40 and 10. The MPLS packet flow #2 is received via fiber optic
 8 cable 115 and passes through optics module 119, serdes 121, framer 123 and
 9 classification engine 124. The MPLS packet is received onto ingress MS-SAR 125
 10 in 16-bit pieces via SPI-4 bus 135 via terminals 198 (see Fig. 10). The MPLS packet
 11 includes a four-byte MPLS header and a data payload. To illustrate operation of the
 12 ingress MS-SAR, the MPLS packet data payload in this example is 150 bytes.

13 Incoming SPI-4 interface block 201 (see Fig. 10) forwards the 154-byte packet
 14 405 as 64-bit pieces to lookup block 202. Figure 40 illustrates the 154-byte packet
 15 405 as it is output from incoming SPI-4 interface block 201. As explained above in
 16 connection with the example of Figure 9, lookup block 202 uses the input port type
 17 supplied via bus 317 to lookup the associated traffic type in the lookup block's "port
 18 table". In this case, the traffic type is "MPLS packet". From the traffic type, lookup
 19 block 202 locates the 20-bit MPLS label within the MPLS header. Lookup block 202
 20 uses the MPLS label to determine the FID and the ingress application type. Lookup
 21 block 202 forwards the FID and ingress application type to the other blocks of the
 22 MS-SAR. The packet 406 is forwarded in 64-bit pieces to segmentation block 203
 23 via bus 318. Although all the 64-bit pieces of a packet for a given input port are
 24 supplied to segmentation block 203 before any other 64-bit pieces are received for
 25 that input port, 64-bit pieces for other input ports may be communicated from lookup
 26 block 202 in an interleaved fashion with the 64-bit pieces of packet 406. Figure 40
 27 illustrates the packet 406 as it is output from lookup block 202.

28 Segmentation block 203 "segments" the incoming packet 406 into 64-byte AAL-5
 29 like chunks. In the present example, there are three such chunks 407-409. The first
 30 chunk 407 contains the 4-byte MPLS header. Segmentation block 203 adds an
 31 AAL5 trailer 410 to the end of the last chunk 409. Any gap between the end of the
 32 26-byte data payload 411 of the last chunk 409 and the trailer 410 is filled with a pad

1 412. The term “segmentation” here does not mean that segmentation block 203 in
 2 this embodiment necessarily breaks the incoming packet into 64-byte chunks. That
 3 has already been done in this embodiment by SPI-4 interface block 201. Rather, the
 4 term “segmentation” means that the packet leaves segmentation block 203 properly
 5 segmented into chunks in accordance with a segmentation scheme. In this
 6 example, segmentation block 203 calculates a cyclic redundancy check (CRC) value
 7 on the chunks of the packet and includes that CRC into trailer 410 of the last chunk
 8 409.

9 Chunks of the packet are output from incoming SPI-4 interface block 201 in 64-
 10 byte chunks. Between the various chunks of a packet, other chunks from other
 11 flows may pass through the lookup block 202 and to the segmentation block 203.
 12 Segmentation block 203 therefore queues the 64-byte chunks 407-409 into a FIFO
 13 on a per port basis. The per port FIFO is located within segmentation block 203. As
 14 segmentation block 203 receives each chunk from an input port, segmentation block
 15 203 reads a partial CRC for the port from a CRC table (not shown) within
 16 segmentation block 203. The partial CRC is modified as a function of the current
 17 chunk received, and the new partial CRC is written back into the CRC table entry for
 18 the appropriate port. When the EOP of a chunk indicates the chunk is the last chunk
 19 of a packet, then the final CRC is calculated and this CRC is inserted into the trailer.

20 Memory manager block 204 reads the 64-byte chunks from the per port FIFO in
 21 segmentation block 203 and stores the 64-byte chunks into buffers in payload
 22 memory 217 using the enqueue command mechanism set forth above in connection
 23 with the example of Figure 9. The per flow queue for flow #2 in this case involves
 24 three buffers.

25 Once the 64-byte chunks for flow #1 and flow #2 are stored in payload memory
 26 217, their per flow queues can be dequeued so as to carry out traffic shaping and/or
 27 output scheduling functions. In the event that the per flow queue for flow #1 is
 28 dequeued, 64-byte chunk 402 (see Figure 39) is output from memory manager block
 29 204 via 128-bit bus 324 (see Figure 10). Reassembly block 205 uses the FID of
 30 chunk 402 to lookup a switch header 413 in the “header table” 327. In this case,
 31 CPU 106 has placed a switch header 413 appropriate for packet-based switch fabric
 32 105 into the location for flow #1 in “header table” 327. Reassembly block 205

removes the 16-byte pad from the first chunk 402, adds the switch packet header 413 to the front to the chunk, and supplies the combination to outgoing SPI-4 interface block 206 using the port calendar mechanism explained above in connection with the example of Figure 9. This outputting to outgoing SPI-4 interface block 206 can occur before all the other 64-byte chunks of the per flow queue are dequeued. In one embodiment of reassembly block 205, either the second 64-byte chunk of a per flow queue or an EOP chunk must be received before the first chunk can be output. The 64-byte chunks 403 and 404 are dequeued and processed by reassembly block 205 in similar fashion, except that no header is added to these chunks. The combination of the switch header 413 and the data portions of the three chunks 402-404 is a switch packet 414.

Figure 39 shows this switch packet 414 as it is output from reassembly block 205 via 64-bit bus 325. The data portions of the three 64-byte chunks 402-404 together form the data payload of switch packet 414. Figure 39 also illustrates switch packet 414 as it is output from outgoing SPI-4 interface block 206.

The 64-byte chunks 407-409 (see Figure 40) of flow #2 are dequeued in similar fashion, the three chunks being processed through the per port queue in reassembly block 205 for the output port of this flow. Reassembly block 205 uses the FID of the first chunk to lookup a switch header 415 in "header table" 327. Reassembly block 205 adds switch header 415 to the front to the first 64-byte chunk. Reassembly block 205 calculates a CRC from the data payload of the three chunks 407-409 and checks to make sure that it matches the CRC found in the trailer 410 of the last chunk 409. After checking the CRC, reassembly block 205 removes the pad 412 and trailer 410.

First chunk 407 is processed and supplied to the outgoing SPI-4 interface block 206 along with switch header 415 before the remaining chunks 408 and 409 of the packet are dequeued. Reassembly block 205 does not store the entire packet in reassembled form on the MS-SAR, but rather processes and outputs the individual 64-byte chunks of the packet one at a time. In this embodiment, the OUTPUT_DATA_FIFO into which reassembly block 205 pushes processed chunks is only 256 bytes in size, a size inadequate to store the entire 160-byte switch packet.

1 Chunks of flow #2 are output from reassembly block 205 using the port calendar
 2 mechanism explained above in connection with the example of Figure 9. The
 3 combination of the data portions of the three chunks 407-409 and the added switch
 4 header 415 together form switch packet 416. Figure 40 shows switch packet 416 as
 5 it is output from reassembly block 205 via 64-bit bus 325. Switch packet 416 is
 6 output from outgoing SPI-4 interface block 206 to the packet switch fabric 105 in 16-
 7 bit pieces via terminals 326.

8 Switch packets 414 and 416 for flow #1 and for flow #2 are switched by the
 9 packet switch fabric 105 of router 100 such that they are supplied to line card 103.
 10 As described above in connection with the example of Figure 9, router 100 uses a
 11 "virtual output port" number in the switch header to identify the particular destination
 12 line card to which the packets should be routed.

13 Processing by egress MS-SAR 200 on flow #1 is explained in connection with
 14 Figures 41, 5 and 10. Switch packet 414 from packet-based switch fabric 105 is
 15 received in 16-bit pieces onto SPI-4 bus 100 via terminals 198. Switch packet 414
 16 includes 8 or 16 bytes of switch header in addition to the packet itself. Figure 41
 17 illustrates switch packet 414 as it is output from incoming SPI-4 interface block 201.
 18 Switch packet 414 is passed to lookup block 202 along with a PID indicating the
 19 logical input port.

20 In the example of Figure 38, flow #1 through egress MS-SAR 200 is a flow from a
 21 packet-based switch fabric that is output as an MPLS packet. As indicated by the
 22 flow to the bottom right of Figure 8, this is egress application type 14. The Azanda
 23 header (see Figure 33) is the last four bytes of the switch header. This Azanda
 24 header contains both the FID as well as the egress application type to be performed
 25 on flow #1 by egress MS-SAR 200. Again, as explained above in connection with
 26 Figure 8, only one traffic type is permitted on any one logical input port. Lookup
 27 block 202 uses the logical input port number supplied to it via 64-bit bus 317 (see
 28 Figure 10) to locate the Azanda header in the information coming from the switch
 29 fabric. Lookup block 202 extracts the FID (FID1) and egress application type
 30 (egress application type 14). Lookup block 202 removes the switch header 413 and
 31 supplies the data payload of the packet via 64-bit bus 318 to segmentation block 203
 32 along with the extracted FID and egress application type.

Figure 41 illustrates the data payload 417 as it is output from lookup block 202. Payload 417, however, is passed to segmentation block 203 in 64-byte chunks. In the example of Figure 38, the 144-byte payload of flow #1 contains enough information for there to be three such chunks. The first chunk 418 contains 64 bytes, the second chunk 419 contains 64 bytes, and the third chunk 420 contains sixteen bytes of data 421. Segmentation block 203 calculates a CRC on the data of the three chunks and adds a trailer 422 so that the trailer is located at the end of the third 64-byte chunk 420. Segmentation block 203 adds any necessary pad between the end of the data 421 of the last chunk and the trailer 422 of the last chunk. The three chunks 418-420 of Figure 41 are queued on a per port basis into an SRAM in segmentation block 203 as explained above in connection with Figure 40 and pointers to the chunks are pushed onto a FIFO in segmentation block 203. Memory manager block 204 pops the FIFO, reads the chunks, and stores the chunks into payload memory 217. Figure 41 illustrates the three 64-byte chunks 418-420 supplied by segmentation block 203 to memory manager block 204. The enqueue command mechanism as outlined above is used such that a per flow queue of BIDs is formed for the chunks of flow #1.

Processing by egress MS-SAR 200 on flow #2 is explained in connection with Figures 42, 5 and 10. The switch packet 416 from packet-based switch fabric 105 is received in 16-bit pieces onto SPI-4 bus 100 via terminals 198. Switch packet 416 includes 8 or 16 bytes of switch header 415, 4 bytes of MPLS header, and 150 bytes of MPLS data payload.

Figure 42 illustrates switch packet 416 as it is output from incoming SPI-4 interface block 201. Switch packet 416 is passed to lookup block 202 along with the logical input port. The lookup block 202 uses the logical input port supplied to it via 64-bit bus 317 to locate the Azanda header in the information coming from the switch fabric. The lookup block extracts the FID and egress application type (egress application type 13). In the case of flow #2, information from a packet-based switch fabric is output from egress MS-SAR 200 as AAL5 cells. This corresponds to egress application type 13 as indicated in the middle right portion of Figure 8. Lookup block 202 removes switch header 415 and supplies the remainder of the packet (MPLS header and MPLS payload) to segmentation block 203 along with the extracted FID

1 and egress application type. Figure 41 illustrates the MPLS header 423 and MPLS
2 data payload 424 as supplied by lookup block 202 to segmentation block 203 via 64-
3 bit bus 318.

4 In egress application type 13, segmentation block 203 “segments” the MPLS
5 header 423 and MPLS data payload 424 into four 48-byte chunks 425-428, because
6 each of these four chunks will be the data payload of an AAL5 cell to be output from
7 egress MS-SAR 200. In this example, the fourth 64-byte chunk contains only ten
8 bytes of data 429. Segmentation block 203 adds a trailer 430 so that the trailer 430
9 is located at the end of the first 48-bytes of the fourth 64-byte chunk. Any necessary
10 pad 431 is added between the end of data 429 and the start of trailer 430.

11 Segmentation block 203 adds a 16-byte pad to the end of each 48-byte chunk to pad
12 each chunk up to the 64-byte size stored in payload memory 217. Figure 42
13 illustrates the four 64-byte chunks as output from segmentation block 203. The
14 mechanism for supplying the four chunks to memory block 204 involves a queue as
15 explained above. Memory manager block 204 pops the queue, receives the 64-byte
16 chunks from segmentation block 203, and stores the 64-byte chunks into 64-byte
17 buffers using the enqueue command mechanism described above. A per flow
18 queue for flow #2 is created.

19 The pointers in the per flow queues for flow #1 and flow #2 are then popped off in
20 a desired order using the dequeue command. Here in this example of Figure 38, we
21 will consider the per flow queue of flow #1 being dequeued first. The three 64-byte
22 chunks for flow #1 are received from payload memory 217 by reassembly block 205
23 and are pushed onto a queue in reassembly block 205 for the intended output port.
24 The pad and trailer of the third 64-byte chunk 420 (see Fig. 41) are removed.
25 Reassembly block 205 also computes a CRC on the combined data of the three
26 chunks and verifies that the CRC in trailer 422 matches the newly computed CRC.
27 For the first 64-byte chunk 418, reassembly block 205 uses the FID to lookup an
28 MPLS header 432 in header table 327 (see Figure 10). The MPLS header 432 was
29 placed in header table 327 by CPU 106 before the lookup operation. Reassembly
30 block 205 performs an MPLS manipulation operation if controlled to do so by control
31 information in the header table.

Figure 41 illustrates a simplified view of MPLS header 432 and the aggregated 144 bytes of data 433 as it is output from reassembly block 205. MPLS header 432 and data 433 together form an MPLS packet 434. Although shown here reassembled as a 144-byte block, the various 64-byte chunks of packet 434 are output from reassembly block 205 one by one, the first 64-byte chunk having the MPLS header appended to it. The chunks are output from reassembly block 205 using the port calendar and queue mechanism set forth above. Outgoing SPI-4 interface block 206 pops the Q_FIFO of reassembly block 205, receives the chunks of the MPLS packet 434 in 64-bit pieces via 64-bit bus 325, and supplies the pieces of the MPLS packet onto terminals 326. The pieces of MPLS packet 434 pass over SPI-4 bus 141 to framer 142, through framer 142, serdes 143, and optics module 144, and onto fiber optic cable 118. It is therefore seen that ATM cell information received onto the line card 101 passes through the packet-based switch fabric and exits line card 103 as an assembled MPLS packet.

The per flow queue for flow #2 is dequeued next in this example (see Figure 42), the 64-byte chunks 425-428 of flow #2 being supplied to the appropriate output port queue in reassembly block 205. In egress application type 13, reassembly block 205 removes the 16-bytes of pad from each 64-byte chunk to recover the 48-bytes of data. Reassembly block 205 uses the FID of a chunk to lookup a 4-byte ATM header 435, adds the ATM header 435 to the 48-byte chunk of data, and adds a 4-byte alignment pad such that each chunk is 56-bytes. The data 429, pad 431, and trailer 430 of the fourth 64-byte chunk 428 form the data portion of the last 56-byte chunk as illustrated in Figure 42. The ATM header 435 is the same for each of the four 56-byte chunks of this flow.

Reassembly block 205 performs ATM translation by using the FID to access the FID entry in "header table" 327 (see Fig. 10). Because the header type in CONTROL_WORD_1 of the FID entry is "ATM", reassembly block 206 replaces the PTI, CLP, EFCI and OAM bits in the ATM header using the incoming EOP, CLP, EFCI and OAM bits from memory manager block 204. The last AAL5 cell (marked EOP) is marked using the PTI field in the ATM header. The packet length is checked against the maximum MTU. Reassembly block 205 checks the length and CRC in AAL5 trailer 430. Figure 42 illustrates the AAL5 cells 436-439 as output

1 from reassembly block 205 in 64-bit pieces onto 64-bit bus 325. Outgoing SPI-4
2 interface block 206 pops the queue for the correct output port, retrieves each AAL5
3 cell, removes the 4-byte alignment pad, and outputs the resulting 52-byte AAL5 cells
4 onto terminals 326 in 16-bit pieces. The AAL5 cells of flow #2 pass over SPI-4 bus
5 141, through framer 142, through serdes 143, through optics module 144, and onto
6 fiber optic cable 118. It is therefore seen that MPLS packet information received
7 onto the line card 101 passes through the packet-based switch fabric and exits line
8 card 103 in AAL5 cell form.

10 APPLICATION TYPES 1 AND 9:

11 Figure 43 illustrates an example of ingress application type 1 and egress
12 application type 9. The incoming data in this example is two AAL5 cells. These two
13 52-byte AAL5 cells are received onto ingress MS-SAR 125 via terminals 198. The
14 two cells are output from the incoming SPI-4 interface block. The first cell contains 4
15 bytes of ATM header and 48 bytes of ATM data #1. The first part of the ATM data
16 #1 in this example is a packet header (for example, an IP header). The second cell
17 contains 4 bytes of ATM header and a small amount of data (ATM data #2). An
18 AAL5 trailer is disposed at the end of the 52 bytes of the second cell, and a pad is
19 disposed between the end of the AAL5 data #2 and the beginning of the trailer in
20 accordance with the AAL5 protocol.

21 Incoming SPI-4 interface block 201 adds four bytes of pad to each cell prior to
22 sending the cells to lookup block 202. The 4-byte pad is added because the data
23 path from the incoming SPI-4 interface block to the segmentation block is 64-bits
24 wide. Segmentation block 203 removes the 4-byte ATM header and adds an
25 additional twelve bytes of pad for a total of sixteen bytes of pad. Memory manager
26 block 204 passes the resulting two 64-byte chunks to payload memory 217 for
27 storage. The two 64-byte chunks are read from payload memory and are passed to
28 reassembly block 205. Reassembly block 205 does not in this ingress application
29 check the CRC in the trailer of the second cell. Reassembly block 205 adds either 8
30 or 16 bytes of a switch header to each 64-byte chunk, and passes the resulting
31 chunks to outgoing SPI-4 interface block 206 as illustrated. Chunks are, however,
32 processed through reassembly block 205 one at a time. As illustrated, the packet

header in the AAL5 data #1 is passed through the various blocks of the ingress MS-SAR to the outgoing SPI-4 interface block 206. The cells pass out of the ingress MS-SAR 125 as “switch cells”, pass through the cell-based switch fabric, and to a line card having an egress MS-SAR. The arrow from the top portion of Figure 43 to the bottom portion of Figure 43 illustrates this passing through the switch fabric.

The two switch cells pass through the incoming SPI-4 interface block and to lookup block 202 of the egress MS-SAR 200 as illustrated. Lookup block 202 removes the switch headers. The resulting two chunks pass through segmentation block 203 and are passed to memory manager block 204 one at a time for storage. Segmentation block 203 in egress application type 9 does not perform segmentation. The two 64-byte chunks are read out of memory manager block 204 and pass to reassembly block 205. Reassembly block 205 maps AAL5 data #1 and AAL5 data #2 into one contiguous data portion. As each of the data portions passes into reassembly block 205, a CRC on the data is computed. The composite CRC is then compared with the CRC in the trailer of the last 64-byte chunk (in this case the second chunk) as that second chunk passes into reassembly block 205. Reassembly block 205 adds an MPLS header and if programmed to do so by control information in the header table, performs an MPLS manipulation operation on the MPLS header. The resulting packet (MPLS header and AAL5 data #1 and AAL5 data #2) is then sent out to the framer via the outgoing SPI-4 interface block 206.

APPLICATION TYPES 2 AND 10:

Figure 44 illustrates an example of ingress application type 2 and egress application type 10. In this example an incoming MPLS packet (MPLS header and 72 bytes of MPLS payload) is received onto the ingress MS-SAR 125 and passes through incoming SPI-4 interface block 201. Lookup block 202 performs a lookup and determines that the ingress application type is ingress application type 2. Segmentation block 203 outputs the packet as two 64-byte chunks such that the first 48 bytes of the first 64-byte chunk includes the MPLS header and a first part (AAL5 data #1) of the data payload. The remainder of the data of the MPLS packet (AAL5 data #2) is segmented into the second 64-byte chunk. Segmentation block 203 calculates a CRC on the 72 bytes of packet data and includes an AAL5 trailer into

1 the second 64-byte chunk as illustrated so that the trailer is located at the end of the
2 first 48 bytes of the second chunk. Segmentation block 203 appends 16 bytes of
3 pad to the end of each 48-byte chunk so that each chunk is padded up to 64 bytes.
4 The resulting two 64-byte chunks pass out of segmentation block 203 one at a time
5 and are stored in memory manager block 204. The 64-bytes chunks are read out of
6 payload memory and pass to reassembly block 205 one at a time. Reassembly
7 block 205 performs a switch header lookup and adds a switch header to each 64-
8 byte chunk. Reassembly block 205 does not perform any checking of the CRC in
9 the trailer. The resulting switch cells are output via outgoing SPI-4 interface block
10 206 as illustrated.

11 The switch cells pass through a switch-based fabric and to an egress line card
12 that includes the egress MS-SAR 200. The switch cells pass through the incoming
13 SPI-4 interface block 201 of the egress MS-SAR 200. Lookup block 202 locates the
14 Azanda header at the end of the switch header and determines the egress
15 application type (egress application type 10) to be performed. Lookup block 202
16 removes the switch headers and passes the remaining two 64-byte chunks to the
17 memory manager block 204 for storage into payload memory. The 64-byte chunks
18 are read out of payload memory and are passed one at a time to reassembly block
19 205. Reassembly block 205 performs a header lookup and adds an ATM header to
20 each 64-byte chunk. Reassembly block 205 removes 12 bytes of pad from the end
21 of each chunk as illustrated in Figure 44. The resulting chunks are passed to
22 outgoing SPI-4 interface block 206. Outgoing SPI-4 interface block 206 last the
23 remaining 4-bytes of pad from the end of each chunk and outputs the chunks as
24 ATM cells to the framer of the egress line card.

25 APPLICATION TYPES 4 AND 14:

26 Figure 45 illustrates an example of ingress application type 4 and egress
27 application type 14. In this example an incoming ATM cell (four bytes of ATM
28 header and 48 bytes of ATM payload) is received onto the ingress MS-SAR 125 and
29 passes through incoming SPI-4 interface block 201. Lookup block 202 performs a
30 lookup and determines that the ingress application type is ingress application type 4.
31 Lookup block 202 adds four bytes of pad as illustrated in Figure 45 and passes the
32

1 resulting 56 bytes to segmentation block 203. Segmentation block 203 generates
2 neither a CRC nor a trailer. Segmentation block 203 adds an additional 8 bytes of
3 pad as illustrated in Figure 45. The resulting 64-byte chunk is stored by memory
4 manager block 204 in payload memory. The 64-byte chunk is read out of payload
5 memory and passes to reassembly block 205. Reassembly block 205 performs a
6 header lookup and adds an 8 or 16-byte switch header to the 64-byte chunk. The
7 resulting switch packet (8 or 16 bytes of switch header, four bytes of ATM header,
8 48 bytes of ATM data, and 12 bytes of pad) is passed to a packet-based switch
9 fabric via outgoing SPI-4 interface block 206. The switch packet passes through the
10 switch fabric and is received onto the appropriate line card and the ingress MS-SAR
11 200. Lookup block 202 locates the Azanda header and determines the egress
12 application type (egress application type 14). Lookup block 202 removes the switch
13 header and passes the remaining 64-bytes to segmentation block 203.
14 Segmentation block 203 does not attach a trailer, but merely passes the 64-byte
15 chunk to memory manager block 204 for storage in payload memory. The 64-byte
16 chunk is retrieved from payload memory and is passed to reassembly block 205.
17 Reassembly block 205 performs a header lookup, and appends an MPLS header as
18 illustrated. Reassembly block 205 removes the last eight bytes of pad. The
19 resulting "encapsulated" ATM cell (MPLS header, 4 bytes of ATM header, 48 bytes
20 of ATM data, and 4 bytes of pad) is output as an MPLS packet to the framer of the
21 egress line card via outgoing SPI-4 interface block 206.

22 APPLICATION TYPES 6 AND 12:

23 Figure 46 illustrates an example of ingress application type 6 and egress
24 application type 12. In this example an incoming ATM encapsulated cell (four bytes
25 of MPLS header, 4 bytes of ATM header, and 48 bytes of ATM payload) is received
26 onto the ingress MS-SAR 125 and passes through incoming SPI-4 interface block
27 201. Lookup block 202 performs a lookup and determines that the ingress
28 application type is ingress application type 6. Lookup block 202 adds four bytes of
29 pad as illustrated in Figure 46 and passes the resulting 56-byte chunk to
30 segmentation block 203. Segmentation block 203 pads the chunk up to 64-bytes but
31 generates neither a CRC nor a trailer. The 64-byte chunk is stored by memory
32

1 manager block 204 in payload memory. The 64-byte chunk is read out of payload
 2 memory and passes to reassembly block 205. Reassembly block 205 performs a
 3 header lookup, adds an 8 or 16-byte switch header to the 64-byte chunk, and
 4 removes the last 8 bytes of pad. The resulting packet (8 or 16 bytes of switch
 5 header, four bytes of MPLS header, four bytes of ATM header, 48 bytes of ATM
 6 data, and 4 bytes of pad) is passed to a packet-based switch fabric via outgoing SPI-
 7 4 interface block 206. The switch packet passes through the packet-based switch
 8 fabric and is received onto the appropriate line card and the egress MS-SAR 200. In
 9 the egress MS-SAR, the lookup block 202 locates the Azanda header and
 10 determines the egress application type (egress application type 12). Lookup block
 11 202 removes the switch header and MPLS header and passes the remaining 56-
 12 bytes (4 bytes of ATM header, 48 bytes of ATM data, and 4 bytes of pad) to
 13 segmentation block 203. Segmentation block 203 does not attach a trailer, but adds
 14 an additional 8-byte pad to pad the chunk up to 64 bytes. The 64-byte chunk is
 15 stored in payload memory, is read out of payload memory, and passes to
 16 reassembly block 205. Reassembly block 205 translates the ATM header. For
 17 example, if instructed so by control information in the header table, reassembly block
 18 205 replaces the VPI/VCI information in the ATM header with VPI/VCI information
 19 retrieved from the header table. Reassembly block 205 removes the last 8 bytes of
 20 pad and passes the resulting ATM cell (4 bytes of ATM header, 48 bytes of data,
 21 and 4 bytes of pad) to outgoing SPI-4 interface block 206. Outgoing SPI-4 interface
 22 block 206 removes the last 4 bytes of pad and passes the resulting 52-byte ATM cell
 23 to the framer of the egress line card.

25 APPLICATION TYPES 6 AND 14:

26 Figure 47 illustrates an example of ingress application type 6 and egress
 27 application type 14. In this example an incoming MPLS packet (four bytes of MPLS
 28 header and 160 bytes of packet payload) is received onto the ingress MS-SAR 125
 29 of a first line card and passes through incoming SPI-4 interface block 201. Lookup
 30 block 202 performs a lookup operation and determines that the ingress application
 31 type is ingress application type 6. Segmentation block 203 outputs the packet
 32 information as 64-byte chunks as illustrated in Figure 47. In the example of Figure

47, the 160 bytes of data payload of the packet is segmented so that the 4-byte MPLS header and 60 bytes of data is segmented into the first 64-byte chunk. The next 64 bytes of data is segmented into the next 64-byte chunk. The remaining 36 bytes of data 440 from the 160-byte payload is segmented into the third 64-byte chunk as illustrated. Segmentation block 203 calculates a CRC on the data payload and places the CRC calculated into a trailer 441 that is incorporated as the end of the last 64-byte chunk. Any space between the end of the data 440 and trailer 441 is filled with a pad 442. Segmentation block 203 may be outputting 64-byte chunks of a packet at the same time that it is receiving chunks of the same packet. The 64-byte chunks are stored by memory manager block 204 into payload memory. The 64-byte chunks are read out of payload memory and are supplied to reassembly block 205. Reassembly block 205 maps the data portions of the 64-bytes chunks into 160-bytes of packet data as illustrated. Reassembly block 205 calculates a CRC on the data as the various chunks of data pass into the reassembly block 205. When the data from the last 64-byte chunk is received onto reassembly block 205, reassembly block compared the final CRC with the CRC in the trailer of the last 64-byte chunk. Using the FID of the flow, reassembly block 205 performs a header lookup and attaches a switch header to the front of the data payload. The resulting switch packet (switch header, MPLS header, and packet data) is supplied to outgoing SPI-4 interface block 206 in pieces, one at a time. Although the packet data payload is illustrated in Figure 46 as being reassembled by reassembly block 205, a first part of the switch packet may be flowing out of reassembly block 205 to outgoing SPI-4 interface block 206 at the same time that a second part of the data portion of the switch packet is being received onto reassembly block 205.

The switch packet of Figure 46 passes through a packet-based switch fabric to a second line card having an egress MS-SAR. The arrow in Figure 46 illustrates this flow of information. The switch packet (switch header, MPLS header, and 160 bytes of packet data in this example) is received onto the egress MS-SAR via incoming SPI-4 interface block 201. Lookup block 202 removes the switch header, locates the Azanda header, and determines from the Azanda header the egress application type to be performed (egress application type 14). Segmentation block 203 segments the packet as illustrated in Figure 46 into 64-byte chunks. Segmentation block adds a

1 trailer to the end of the last 64-byte chunk as explained above. The 64-byte chunks
2 are supplied one by one to the memory manager block 204 and are stored one by
3 one into payload memory. The 64-bytes are read out of payload memory one by
4 one and are supplied, one by one, to reassembly block 205. When reassembly
5 block 205 receives the first 64-byte chunk, reassembly block 205 performs a header
6 lookup using header table 327. If instructed to by control information found in the
7 header table 327, reassembly block 205 performs an MPLS manipulation operation
8 on the MPLS header. For example, reassembly block 205 may replace the MPLS
9 label in the MPLS header with an MPLS retrieved from header table 327.

10 Reassembly block 205 maps the data portions of the three 64-byte chunks into one
11 packet payload as illustrated in Figure 46. As the data portion of the various 64-byte
12 chunks pass through reassembly block 205, reassembly block 205 calculates a CRC
13 on the data such that reassembly block 205 checks the final composite CRC with the
14 CRC found in the last 64-byte chunk. The first part of the resulting MPLS packet can
15 be output from reassembly block 205 at the same time that later 64-byte chunks for
16 the packet are being received onto reassembly block 205. The resulting MPLS
17 packet (MPLS header, and 160 bytes of packet data) is output to the framer of the
18 egress line card via outgoing SPI-4 interface block 206.

19 20 REASSEMBLY "ON-THE-FLY":

21 Where a segmented packet is received on an egress line card in the form of a
22 plurality of individual cells (for example AAL5 cells), and where the data payloads of
23 the individual cells are to be reassembled into the original packet, the data payloads
24 of the individual cells are buffered in payload memory before any reassembly is
25 done. Rather than doing reassembly before the cells are buffered, reassembly in
26 the MS-SAR 200 is done on a per output port basis as the cells for a flow are being
27 output (i.e., on-the-fly) from the egress MS-SAR. A packet being reassembled is not
28 necessarily reassembled in the sense that the entire reassembled packet exists in
29 reassembled form in reassembly block 205. Rather the individual 64-byte chunks
30 making up the packet are processed by reassembly block 205 one at a time. It is,
31 for example, possible that a first part of the packet has already passed out of the

1 reassembly block 205 before the last part of the packet has been read out of
2 payload memory 204.

3 In accordance with one embodiment, a "reassembly context" is maintained for
4 each reassembly process going on. A reassembly context may, for example,
5 include a 32-bit partial CRC and 16-bit partial packet byte count. Rather than
6 maintaining such a reassembly context for each of the many flows coming into the
7 egress MS-SAR (it is, for example, possible that each incoming flow would be AAL5
8 cells requiring reassembly), reassembly is done on a per output port basis after
9 buffering in payload memory 217, thereby reducing the maximum number of
10 reassembly contexts that have to be maintained to one per output port. There is at
11 most one such packet reassembly going on for each of the output ports that is
12 configured and active. Reducing the number of reassembly contexts to be stored
13 reduces the amount of memory necessary and thereby reduces line card costs.

14 15 SEGMENTATION "ON-THE-FLY":

16 A packet received onto an ingress line card may, in accordance with an
17 adaptation layer protocol, require segmentation into a plurality of pieces, where each
18 of the pieces is the data payload of an individual cell. The AAL5 protocol is an
19 example of one such protocol. The segmenting is done so that the individual cells
20 can later be reassembled into the original packet. Rather than storing such a packet
21 into payload memory, retrieving the packet, and then segmenting it as it is sent out
22 to the switch fabric, segmentation in ingress MS-SAR 125 is done on a per input port
23 basis as flows are received (i.e., on-the-fly) onto the ingress MS-SAR.

24 Segmentation, as the term is used here, encompasses the one-by-one processing of
25 segments of a packet such that the processed segments include information
26 required by the segmentation protocol. A segmentation block may, for example,
27 output a processed segment before the last part of the packet has even been
28 received by the segmentation block.

29 In accordance with one embodiment, a segmentation context is maintained for
30 each segmentation process going on. A segmentation context may, for example,
31 include a 32-bit CRC value and a 16-bit packet byte count value. Because a packet
32 to be segmented is received onto segmentation block 203 in segments, and

1 because segmentation block 203 processes these segments one by one, the
2 segmentation engine maintains a partial 32-bit CRC value for the packet and a
3 partial packet byte count value. As each segment is processed, segmentation block
4 203 updates the partial CRC value and partial packet byte count value. After the last
5 segment has been received onto the segmentation block as indicated by the EOP bit
6 of the segment, then the partial CRC value is the complete CRC value for the
7 packet. Similarly, the partial packet byte count value contains the number of bytes in
8 the packet. The 32-bit complete CRC value and the 16-bit packet byte count value
9 are included by segmentation block 203 into the AAL5-like trailer that the
10 segmentation block appends to the last 64-byte chunk for the packet. Rather than
11 maintaining such a segmentation context for each of the many flows going out of the
12 ingress MS-SAR (it is, for example, possible that each outgoing flow is to be in the
13 form of AAL5 cells), this segmentation process is done before buffering on a per
14 input port basis, thereby reducing the maximum number of segmentation contexts
15 that have to be maintained to one per input port. There is at most one such packet
16 segmentation process going on for each of the input ports that is configured and
17 active. Reducing the number of segmentation contexts to be stored reduces the
18 amount of memory necessary and thereby reduces line card costs. If a separate
19 counter or separate CRC engine is provided for each segmentation process, then
20 reducing the number of segmentation processes going on at one time further
21 reduces costs.

23 MULTIPLE DATA PATH CHIPS TO INCREASE THROUGHPUT RATE:

24 Figure 48 illustrates a novel aspect wherein MS-SAR functionality is partitioned
25 into a control integrated circuit and a data path integrated circuit such that system
26 throughput can be increased by using multiple data path integrated circuits. This
27 increase in system throughput is accomplished without having to redesign either the
28 data path integrated circuit or the control integrated circuit. The data path integrated
29 circuit and the control integrated circuit therefore make a versatile chip set. Where a
30 lower throughput is required, cost is reduced by using just one data path integrated
31 circuit with one control integrated circuit. Where a higher throughput is required,
32 multiple data path integrated circuits are used with one control integrated circuit.

1 Accordingly, the very same data path and control integrated circuits are useful for a
2 wide range of throughput applications. As a result, larger production runs of each
3 integrated circuit is likely and an associated reduction in per part cost is expected.

4 Figure 48 illustrates a line card 500 having a first MS-SAR 501 configured in the
5 ingress mode and a second MS-SAR 502 configured in the egress mode. Network
6 information passing from fiber optic cable 503 to the switch fabric passes through
7 optics module 504, serdes 505, framer/mapper 506 and classification engine 507
8 before it reaches ingress MS-SAR 501. In an egress path, network information
9 passes from egress MS-SAR 502, through framer/mapper 506, serdes 508 and
10 optics module 509 to fiber optic cable 510. The optics/serdes/framer circuitry of
11 Figure 48 is similar to the corresponding circuitry of Figure 5, except that the circuitry
12 in Figure 48 can receive from and output to fiber optic cables 503 and 510 are higher
13 data throughput rates. For example, where line card 101 of Figure 5 may receive
14 and output at OC-192 line rates (about 10 gigabits per second), line card 500 of
15 Figure 48 may receive and output at OC-768 line rates (about 40 gigabits per
16 second).

17 How the partitioning of MS-SAR functionality facilitates handling higher data
18 throughput rates is explained in connection with ingress MS-SAR 501 of Figure 48.
19 Ingress MS-SAR 501 includes a distribution integrated circuit 511, a control
20 integrated circuit 512, four data path integrated circuits 513-516, and one
21 aggregation integrated circuit 517. Control integrated circuit 512 includes the
22 circuitry shown in Figure 10 within dashed line 213. Each of the data path integrated
23 circuits 513-516 includes the circuitry shown in Figure 10 within dashed line 212.
24 Although per flow queue block 207 in the embodiment of Figure 10 has only one
25 control interface for sending and receiving enqueue and dequeue commands, the
26 control integrated circuit 512 of the embodiment of Figure 48 has four such
27 interfaces. Accordingly, ingress control integrated circuit 512 of Figure 48 is coupled
28 to data path integrated circuits 513-516 via control buses 518-521, respectively.

29 Operation of the circuit of Figure 48 is explained in connection with an example of
30 a high-speed stream of multiple packets being received from fiber optic cable 503 at
31 OC-768 rates. The packets pass, one by one, through optics module 504, serdes
32 505, framer/mapper 506, and classification engine 507 to distribution integrated

1 circuit 511. Distribution integrated circuit 511 receives the packets from
2 classification engine 507 via an SPI-4-like bus 522. This bus 522 is like an SPI-4
3 bus, except that it is capable of operating at 40 gigabits per second rates. Each
4 packet is framed by a start delimiter and an end delimiter.

5 Distribution integrated circuit 511 distributes each incoming packet to one of the
6 four data path integrated circuits 513-516, the particular data path integrated circuit
7 being chosen so as to distribute data path processing load evenly among the
8 multiple data path integrated circuits. Any of many known load-balancing algorithms
9 can be employed. Distribution integrated circuit 511 only outputs complete packets.
10 In addition to distributing load, distribution integrated circuit 511 includes a sequence
11 number along with each packet. In one embodiment, there is a sequence of
12 sequence numbers for each flow. Distribution integrated circuit 511 employs a
13 lookup block like lookup block 202 of Figure 10. Lookup information is provisioned
14 by the CPU beforehand as explained above in connection with lookup block 202 of
15 Figure 10. Distribution integrated circuit 511 receives a packet, identifies the flow,
16 increments a sequence number for that flow, and adds the sequence number to the
17 packet. Each respective packet of a flow therefore carries a sequence number that
18 is one greater than the sequence number of the previous packet in that flow.

19 Figure 49 is a diagram of a packet as output by distribution integrated circuit 511.
20 The packet includes a data portion 523 and a header portion 524. The packet is
21 framed by a start delimiter 525 and an end delimited 526. The sequence number
22 527 added by distribution integrated circuit 511 is disposed between the start
23 delimited 525 and the packet header 524.

24 The various packets, after being marked with sequence numbers, flow into the
25 various data path integrated circuits 513-516. The lookup block of the data path
26 integrated circuit that receives the packet determines the flow ID of the packet as
27 described above in connection with the embodiment of Figure 10. In the
28 embodiment of Figure 48, the lookup block also extracts the packet sequence
29 number. Both the flow ID and the packet sequence number are then supplied to the
30 control integrated circuit 512 via the enqueue mechanism. Control integrated circuit
31 512 is therefore aware of which packets of which flows have been received onto
32 which data path integrated circuits. In addition to performing the other traffic

1 management, policing, shaping, and metering functions described above, control
2 integrated circuit 512 issues dequeue commands to ensure that the various packets
3 of a flow are supplied to aggregation integrated circuit 517 in the same order in
4 which they were received onto distribution integrated circuit 511. To do this, control
5 integrated circuit 512 maintains a "packet queue" for each flow.

6 Figures 50 and 51 illustrate the building of such a packet queue. In Figure 50,
7 the packets of a flow (FID1) were sequence numbered P1, P2, P3 and so forth by
8 distribution integrated circuit 511. Control integrated circuit 512 maintains a head
9 pointer and a tail pointer for this packet queue. The first packets P1, P2, P3 in this
10 flow have already been dequeued and sent through the aggregation integrated
11 circuit 517 to the switch fabric. Consequently they do not appear in the packet
12 queue. The next packet to be dequeued is packet P5. Control integrated circuit 512
13 will dequeue the various 64-byte chunks of this packet by dequeuing a per flow
14 queue for this packet as explained above in connection with Figure 10. In the
15 example of Figure 50, packet P7 of flow #1 has not been queued in the packet
16 queue. Accordingly, control integrated circuit 512 will not dequeue packet P8, but
17 rather will wait until packet P7 has been queued.

18 Figure 51 illustrates the packet queue after packet P7 has been queued. Note
19 that packet P7 has been linked into the queue in its proper place before packet P8
20 such that the dequeuing of the linked list will result in the packets of the flow being
21 dequeued in the proper order. In the example of Figure 51, packet P5 has been
22 dequeued, and a new packet for this flow, packet P10, has been received and
23 queued.

24 Control integrated circuit 512 dequeues the various packets in accordance with
25 this scheme so that the various packets of each flow are output to aggregation
26 integrated circuit 517 in the correct order. Aggregation integrated circuit 517
27 combines the packets in the order it receives them into one stream. That one
28 stream is output to the switch fabric via SPI-4-like bus 528. Bus 528 is like an SPI-4
29 bus, except that it is capable of operating at 40 gigabits per second rates. It is
30 therefore seen that one high throughput rate data path (40 gigabits/second) coming
31 into the line card is processed using four lower throughput rate data path integrated
32 circuits (10 gigabits/second), and that the outputs of the four lower throughput rate

1 data path integrated circuits are combined to form one higher throughput rate data
2 path (40 gigabits/second) to the switch fabric.

3 In the embodiment of Figure 48, the incoming data path having the increased
4 throughput rate is being controlled by only one control integrated circuit 512. The
5 associated increase in processing required of control integrated circuit 512 may
6 result in difficulties in accessing external memory.

7 Figure 52 is a diagram of an external memory device 529 that is coupled to a
8 control integrated circuit. External memory device 529 may, for example, be
9 external memory 225 of Figure 10. External memory device 529 stores two types of
10 information, information #1 and information #2, both of which must be accessed
11 within a particular amount of time related to the rate of incoming information. Where
12 the control integrated circuit is clocked by a clock signal, this particular amount of
13 time can be referred to as eight clock periods. In the example to the left of Figure
14 52, each piece of information can be accessed in two clock periods. Both pieces of
15 information are stored in the same external memory device requiring one to be
16 accessed before the other. A total of four clock periods is therefore required to
17 access both pieces of information.

18 If the circuit of Figure 48 is now employed to handle OC-768 line rate information
19 coming in from fiber optic cable 503, then this same information (information #1 and
20 information #2) must be accessed in a smaller amount of time. Consider the
21 example where both information #1 and information #2 must now be accessed within
22 two clock periods. A memory may not be available that has fast enough access
23 times to handle the accessing required within this fewer number of clock periods.

24 In accordance with one novel aspect, two external memories are used.
25 Information #1 is stored in the first external memory 530 and information #2 is stored
26 in second external memory 531. The two external memories are accessed at the
27 same time in parallel. If external memories 530 and 531 are the same type of
28 external memory 529, then these external memories have access times of two clock
29 periods and both information #1 and information #2 are accessed within the required
30 two clock periods. It is to be understood, of course, that the eight and two in the
31 example are used only to illustrate the technique of accessing memories in parallel
32 to facilitate handling higher data throughput rates. The technique here is not limited

1 to the numbers in this example. An example of information #1 is cell count
2 information stored in the embodiment of Figure 10 in PFQ STAT memory 225. An
3 example of information #2 is packet count information stored in the embodiment of
4 Figure 10 in PFQ STAT memory 225. These two types of information are, in one
5 embodiment of Figure 48, stored in different external memory devices.

6 Although the operation of the embodiment of Figure 48 is described in
7 connection with a flow of packets, the same process is followed to handle flows of
8 cells. The embodiment of Figure 48 handles all the application types of Figures 7
9 and 8. It is also to be understood that the functionality of the circuit of Figure 48 can
10 be rearranged. For example, the distribution chip and the control integrated circuit
11 may be combined onto one integrated circuit such that only one lookup block is
12 required. Alternatively, the lookup function may be performed on the distribution
13 chip such that results of the lookup are forwarded to the control integrated circuit.
14 Various rearrangements of the functionality of Figure 48 are possible without losing
15 the benefit of the novel partitioning of the MS-SAR.

16 17 BACKPRESSING USING SERIAL BUS:

18 A router 600 involves a first line card 601 and a second line card 601. Each of
19 the first and second line cards involves an MS-SAR operating in the ingress mode
20 and an MS-SAR operating in the egress mode. The egress MS-SAR 603 on the
21 second line card can become endangered of being overloaded if, for example, the
22 ingress MS-SAR 604 on the first line card continues to send network information for
23 a flow to the egress MS-SAR 603 on the second line card, but the egress MS-SAR
24 603 on the second line card is prevented from outputting that information, for
25 example due to congestion at the framer 605. Consequently, more and more of the
26 network information for the flow will end up having to be buffered by the egress MS-
27 SAR 603 of the second line card (buffered in payload memory).

28 In one novel aspect, the ingress and egress MS-SAR devices 604 and 606 of the
29 first line card 601 are linked by a serial bus 607 on the first line card, and the ingress
30 and egress MS-SAR devices 603 and 608 of the second line card 602 are linked by
31 a serial bus 609 on the second line card. If the egress MS-SAR 603 of the second
32 line card is in danger of becoming overloaded, then the egress MS-SAR 603 of the

1 second line card sends an indication of this situation to the ingress MS-SAR 608 of
2 the second line card via the serial bus 609 on the second line card. In one
3 embodiment, it is the per flow queue block of the egress MS-SAR 603 that detects
4 that the size of the free buffer queue has decreased to an undesirably low level. The
5 per flow queue block is therefore coupled to the serial bus 609 as illustrated in
6 Figure 53.

7 The ingress MS-SAR 608 of the second line card receives this indication from
8 serial bus 609 and relays the indication to the first line card 601 by outputting a
9 special status switch cell 611. The special status switch cell 611 is transported
10 across the switch fabric 610 to the egress MS-SAR 606 of the first line card. The
11 egress MS-SAR 606 of the first line card detects the special status switch cell, and
12 relays the indication of the situation to the ingress MS-SAR 604 of the first line card
13 via the serial bus 607 on the first line card. In one embodiment, it is the
14 segmentation block that detects the special status switch cell. The segmentation
15 block is therefore coupled to the serial bus 607 as illustrated in Figure 53.

16 In response to receiving this indication from serial bus 607, the ingress MS-SAR
17 604 on the first line card slows or stops its outputting of the information that is
18 overburdening the egress MS-SAR 603 on the second line card. In one
19 embodiment, the output scheduler is able to slow or stop the outputting of
20 information that is overburdening egress MS-SAR 603. Consequently, the serial bus
21 607 is coupled to the output scheduler block of ingress MS-SAR 604 as illustrated in
22 Figure 53. For additional details on this backpressuring mechanism, see: U.S.
23 Patent Application Serial No. 09/779,381, filed February 7, 2001, by Parruck et al.
24 (the subject matter of which is incorporated herein by reference).

25 Although the present invention is described in connection with certain specific
26 embodiments for instructional purposes, the present invention is not limited thereto.
27 Although the bus protocol by which the incoming bus interface block communicates
28 off-chip and the protocol by which the outgoing bus interface communicates off-chip
29 are the same (SPI-4) in the above-described embodiments, they need not be the
30 same. Only one or neither of the interface blocks need communicate off-chip in
31 accordance with the SPI-4 protocol. The SPI-4 bus is mentioned just as an
32 example. More than one incoming bus interface block can be provided so that the

1 MS-SAR can interface to either of two different types of buses. Although the lookup
2 block is disposed in the data path between the incoming bus interface block and the
3 segmentation block, this need not be the case in all embodiments. The lookup
4 block may, for example, analyze data passing from the incoming bus interface to the
5 segmentation block without acting as a buffer in the data path between the incoming
6 bus interface block and the segmentation block. In some embodiments, the
7 segmentation block does buffer an entire packet before segmenting it, rather than
8 receiving the packet in smaller pieces and then processing those smaller pieces one
9 by one. Similarly, the reassembly block in some embodiments does reassemble an
10 entire packet before sending the reassembled packet out to the outgoing bus
11 interface block. Although in the embodiments described above the incoming
12 network data is stored in payload memory in buffers that all have the same size, this
13 need not be the case. Multiple different buffer sizes are employed in some
14 embodiments. The data of packets could, for example, be stored in one size of
15 buffer whereas the data of ATM cells could be stored in another size of buffer.
16 Although the Azanda header is described here as being embedded in the switch
17 header, this is not the case in all embodiments. In one embodiment, the Azanda
18 header is provided in the switch cell after the payload. The Azanda header is
19 intended to provide a general-purpose vehicle for the communication of information
20 from an ingress MS-SAR to an egress MS-SAR. Information other than FID and
21 egress application type can be included in such an Azanda header. The specific
22 format for the Azanda header provided above is provided to illustrate the general
23 concept. The architectural concepts disclosed are not limited to any particular
24 protocol, but rather are generally applicable to other protocols. Other application
25 types can be supported, where the various functional blocks of the MS-SAR operate
26 on a flow in a particular way determined by the application type of the flow.
27 Accordingly, various modifications, adaptations, and combinations of various
28 features of the described embodiments can be practiced without departing from the
29 scope of the invention as set forth in the claims.